# CSC 2515 (Lecture 3): Linear Classification

#### Raquel Urtasun & Rich Zemel

University of Toronto

Jan 19, 2015

- Classification as regression
- Decision boundary
- Loss functions
- Logistic Regression
- Regularization
- Nearest Neighbors

- We are interested in mapping the input  $\mathbf{x} \in \mathcal{X}$  to a *label*  $t \in \mathcal{Y}$
- In regression typically  $\mathcal{Y}=\Re$
- Now its a category
- Examples?



What digit is this? How can I predict this? What are my input features?

Urtasun & Zemel (UofT)

CSC 2515



Is this a dog? How can I predict this? What are my input features?



what about this one?



Am I going to pass the exam? How can I predict this? What are my input features?



Do I have diabetes? How can I predict this? What are my input features?

- **Generative approach**: Model the distribution of inputs characteristic of the class (Bayes classifier)
  - build a model of  $p(\mathbf{x}|t = k)$  for every class
  - apply Bayes rule to predict the class
- **Discriminative approach:** estimate parameters of decision boundary/class separator directly from labeled examples
  - learn boundary parameters directly (e.g., logistic regression)
  - learn mappings from inputs to classes  $x \rightarrow t$  (e.g., neural nets)

• Can we do this task using what we have learned in the previous lecture?

- Can we do this task using what we have learned in the previous lecture?
- Simple hack: Ignore that the input is categorical!

- Can we do this task using what we have learned in the previous lecture?
- Simple hack: Ignore that the input is categorical!
- Suppose we have a binary problem,  $t \in \{-1, 1\}$

- Can we do this task using what we have learned in the previous lecture?
- Simple hack: Ignore that the input is categorical!
- Suppose we have a binary problem,  $t \in \{-1, 1\}$
- Assuming the standard model used for regression

$$y = f(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$

- Can we do this task using what we have learned in the previous lecture?
- Simple hack: Ignore that the input is categorical!
- Suppose we have a binary problem,  $t \in \{-1, 1\}$
- Assuming the standard model used for regression

$$y = f(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$

• How can we obtain **w**?

- Can we do this task using what we have learned in the previous lecture?
- Simple hack: Ignore that the input is categorical!
- Suppose we have a binary problem,  $t \in \{-1, 1\}$
- Assuming the standard model used for regression

$$y = f(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$

- How can we obtain **w**?
- Use least squares,  $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$ . How is  $\mathbf{X}$  computed? and  $\mathbf{t}$ ?

- Can we do this task using what we have learned in the previous lecture?
- Simple hack: Ignore that the input is categorical!
- Suppose we have a binary problem,  $t \in \{-1, 1\}$
- Assuming the standard model used for regression

$$y = f(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$

- How can we obtain w?
- Use least squares,  $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$ . How is  $\mathbf{X}$  computed? and  $\mathbf{t}$ ?
- Which loss are we minimizing? Does it make sense?

$$\ell_{square}(\mathbf{w}, t) = \frac{1}{N} \sum_{i=1}^{N} (t_i - \mathbf{w}^T \mathbf{x}_i)^2$$

- Can we do this task using what we have learned in the previous lecture?
- Simple hack: Ignore that the input is categorical!
- Suppose we have a binary problem,  $t \in \{-1, 1\}$
- Assuming the standard model used for regression

$$y = f(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$

- How can we obtain w?
- Use least squares,  $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$ . How is  $\mathbf{X}$  computed? and  $\mathbf{t}$ ?
- Which loss are we minimizing? Does it make sense?

$$\ell_{square}(\mathbf{w}, t) = \frac{1}{N} \sum_{i=1}^{N} (t_i - \mathbf{w}^T \mathbf{x}_i)^2$$

• How do I compute a label for a new example? Let's see an example







• Our classifier has the form

$$f(\mathbf{x}, \mathbf{w}) = w_o + \mathbf{w}^T \mathbf{x}$$



• Our classifier has the form

$$f(\mathbf{x}, \mathbf{w}) = w_o + \mathbf{w}^T \mathbf{x}$$

• A reasonable decision rule is

$$y = egin{cases} 1 & ext{if } f(\mathbf{x}, \mathbf{w}) \geq 0 \ -1 & ext{otherwise} \end{cases}$$



• Our classifier has the form

$$f(\mathbf{x}, \mathbf{w}) = w_o + \mathbf{w}^T \mathbf{x}$$

• A reasonable decision rule is

$$y = egin{cases} 1 & ext{if } f(\mathbf{x}, \mathbf{w}) \geq 0 \ -1 & ext{otherwise} \end{cases}$$

• How can I mathematically write this rule?

$$y = \operatorname{sign}(w_0 + \mathbf{w}^T \mathbf{x})$$



• Our classifier has the form

$$f(\mathbf{x}, \mathbf{w}) = w_o + \mathbf{w}^T \mathbf{x}$$

• A reasonable decision rule is

$$y = egin{cases} 1 & ext{if } f(\mathbf{x}, \mathbf{w}) \geq 0 \ -1 & ext{otherwise} \end{cases}$$

• How can I mathematically write this rule?

$$y = \operatorname{sign}(w_0 + \mathbf{w}^T \mathbf{x})$$

• How does this function look like?



• How can I mathematically write this rule?

$$y = \operatorname{sign}(w_0 + \mathbf{w}^T \mathbf{x})$$

• This specifies a linear classifier: it has a linear boundary (hyperplane)

$$w_0 + \mathbf{w}^T \mathbf{x} = 0$$

which separates the space into two "half-spaces"



$$w_0 + \mathbf{w}^T \mathbf{x} = 0$$

which separates the space into two "half-spaces"



$$w_0 + \mathbf{w}^T \mathbf{x} = 0$$

which separates the space into two "half-spaces"

• In 1D this is simply a threshold

## Example in 2D



• The linear classifier has a linear boundary (hyperplane)

$$w_0 + \mathbf{w}^T \mathbf{x} = 0$$

which separates the space into two "half-spaces"

# Example in 2D



• The linear classifier has a linear boundary (hyperplane)

$$w_0 + \mathbf{w}^T \mathbf{x} = 0$$

which separates the space into two "half-spaces"

In 2D this is a line

Urtasun & Zemel (UofT)



$$w_0 + \mathbf{w}^T \mathbf{x} = 0$$

which separates the space into two "half-spaces"



$$w_0 + \mathbf{w}^T \mathbf{x} = 0$$

which separates the space into two "half-spaces"

In 3D this is a plane



$$w_0 + \mathbf{w}^T \mathbf{x} = 0$$

which separates the space into two "half-spaces"

- In 3D this is a plane
- What about higher-dimensional spaces?

## Geometry

 $\mathbf{w}^T \mathbf{x} = 0$  a line passing though the origin and orthogonal to  $\mathbf{w}$  $\mathbf{w}^T \mathbf{x} + w_0 = 0$  shifts it by  $w_0$ 



#### Figure from G. Shakhnarovich

• Learning consists in estimating a "good" decision boundary

- Learning consists in estimating a "good" decision boundary
- $\bullet$  We need to find  $\boldsymbol{w}$  (direction) and  $\mathit{w}_0$  (location) of the boundary

- Learning consists in estimating a "good" decision boundary
- We need to find  $\mathbf{w}$  (direction) and  $w_0$  (location) of the boundary
- What does "good" mean?

- Learning consists in estimating a "good" decision boundary
- We need to find **w** (direction) and  $w_0$  (location) of the boundary
- What does "good" mean?
- Is this boundary good?


## Learning Linear Classifiers

- Learning consists in estimating a "good" decision boundary
- We need to find **w** (direction) and  $w_0$  (location) of the boundary
- What does "good" mean?
- Is this boundary good?



# Learning Linear Classifiers

- Learning consists in estimating a "good" decision boundary
- We need to find **w** (direction) and  $w_0$  (location) of the boundary
- What does "good" mean?
- Is this boundary good?



• We need a criteria that tell us how to select the parameters

# Learning Linear Classifiers

- Learning consists in estimating a "good" decision boundary
- We need to find **w** (direction) and  $w_0$  (location) of the boundary
- What does "good" mean?
- Is this boundary good?



- We need a criteria that tell us how to select the parameters
- Do you know any?

 $\bullet\,$  Classifying using a linear decision boundary reduces the data dimension to 1

$$y(\mathbf{x}) = \operatorname{sign}(w_0 + \mathbf{w}^T \mathbf{x})$$

 $\bullet\,$  Classifying using a linear decision boundary reduces the data dimension to  $1\,$ 

$$y(\mathbf{x}) = \operatorname{sign}(w_0 + \mathbf{w}^T \mathbf{x})$$

• What is the cost of being wrong?

• Classifying using a linear decision boundary reduces the data dimension to 1

$$y(\mathbf{x}) = \operatorname{sign}(w_0 + \mathbf{w}^T \mathbf{x})$$

- What is the cost of being wrong?
- Loss function: L(y, t) is the loss incurred for predicting y when correct answer is t

• Classifying using a linear decision boundary reduces the data dimension to 1

$$y(\mathbf{x}) = \operatorname{sign}(w_0 + \mathbf{w}^T \mathbf{x})$$

- What is the cost of being wrong?
- Loss function: L(y, t) is the loss incurred for predicting y when correct answer is t
- For medical diagnosis: For a diabetes screening test is it better to have false positives or false negatives?

 $\bullet\,$  Classifying using a linear decision boundary reduces the data dimension to  $1\,$ 

$$y(\mathbf{x}) = \operatorname{sign}(w_0 + \mathbf{w}^T \mathbf{x})$$

- What is the cost of being wrong?
- Loss function: L(y, t) is the loss incurred for predicting y when correct answer is t
- For medical diagnosis: For a diabetes screening test is it better to have false positives or false negatives?
- For movie ratings: The "truth" is that Alice thinks E.T. is worthy of a 4. How bad is it to predict a 5? How about a 2?

• A possible loss to minimize is the zero/one loss

$$L(y(\mathbf{x}), t) = \begin{cases} 0 & \text{if } y(\mathbf{x}) = t \\ 1 & \text{if } y(\mathbf{x}) \neq t \end{cases}$$

• A possible loss to minimize is the zero/one loss

$$L(y(\mathbf{x}), t) = \begin{cases} 0 & \text{if } y(\mathbf{x}) = t \\ 1 & \text{if } y(\mathbf{x}) \neq t \end{cases}$$

• Is this minimization easy to do? why?

• Zero/one loss for a classifier

$$L_{0-1}(y(\mathbf{x}), t) = \begin{cases} 0 & \text{if } y(\mathbf{x}) = t \\ 1 & \text{if } y(\mathbf{x}) \neq t \end{cases}$$

• Zero/one loss for a classifier

$$L_{0-1}(y(\mathbf{x}), t) = \begin{cases} 0 & \text{if } y(\mathbf{x}) = t \\ 1 & \text{if } y(\mathbf{x}) \neq t \end{cases}$$

• Asymmetric Binary Loss

$$L_{ABL}(y(\mathbf{x}), t) = \begin{cases} \alpha & \text{if } y(\mathbf{x}) = 1 \land t = 0\\ \beta & \text{if } y(\mathbf{x}) = 0 \land t = 1\\ 0 & \text{if } y(\mathbf{x}) = t \end{cases}$$

• Zero/one loss for a classifier

$$L_{0-1}(y(\mathbf{x}), t) = \begin{cases} 0 & \text{if } y(\mathbf{x}) = t \\ 1 & \text{if } y(\mathbf{x}) \neq t \end{cases}$$

• Asymmetric Binary Loss

$$L_{ABL}(y(\mathbf{x}), t) = \begin{cases} \alpha & \text{if } y(\mathbf{x}) = 1 \land t = 0\\ \beta & \text{if } y(\mathbf{x}) = 0 \land t = 1\\ 0 & \text{if } y(\mathbf{x}) = t \end{cases}$$

• Squared (quadratic) loss

$$L_{squared}(y(\mathbf{x}), t) = (t - y(\mathbf{x}))^2$$

• Zero/one loss for a classifier

$$L_{0-1}(y(\mathbf{x}), t) = \begin{cases} 0 & \text{if } y(\mathbf{x}) = t \\ 1 & \text{if } y(\mathbf{x}) \neq t \end{cases}$$

• Asymmetric Binary Loss

$$L_{ABL}(y(\mathbf{x}), t) = \begin{cases} \alpha & \text{if } y(\mathbf{x}) = 1 \land t = 0\\ \beta & \text{if } y(\mathbf{x}) = 0 \land t = 1\\ 0 & \text{if } y(\mathbf{x}) = t \end{cases}$$

• Squared (quadratic) loss

$$L_{squared}(y(\mathbf{x}), t) = (t - y(\mathbf{x}))^2$$

Absolute Error

$$L_{quadratic}(y(\mathbf{x}),t) = |t - y(\mathbf{x})|$$

• What if the movie predictions are used for rankings? Now the predicted ratings don't matter, just the order that they imply.

- What if the movie predictions are used for rankings? Now the predicted ratings don't matter, just the order that they imply.
- In what order does Alice prefer E.T., Amelie and Titanic?

- What if the movie predictions are used for rankings? Now the predicted ratings don't matter, just the order that they imply.
- In what order does Alice prefer E.T., Amelie and Titanic?
- Possibilities:
  - 0-1 loss on the winner
  - Permutation distance
  - Accuracy of top K movies.

#### Can we always separate the classes?



#### Can we always separate the classes?



#### Can we always separate the classes?



How can we obtain a non-linear decision boundary?

Causes of non perfect separation

- Model is too simple
- Noise in the inputs (i.e., data attributes)
- Simple features that do not account for all variations
- Errors in data targets (miss labelings)

Causes of non perfect separation

- Model is too simple
- Noise in the inputs (i.e., data attributes)
- Simple features that do not account for all variations
- Errors in data targets (miss labelings)

Should we make the model complex enough to have perfect separation in the training data?

• The classifier we have looked at is

$$y(\mathbf{x}) = \operatorname{sign}(w_0 + \mathbf{w}^T \mathbf{x})$$

• The classifier we have looked at is

$$y(\mathbf{x}) = \operatorname{sign}(w_0 + \mathbf{w}^T \mathbf{x})$$

• It was difficult to optimize any loss on  $\ell(y, t)$  due to the form of  $y(\mathbf{x})$ 

• The classifier we have looked at is

$$y(\mathbf{x}) = \operatorname{sign}(w_0 + \mathbf{w}^T \mathbf{x})$$

- It was difficult to optimize any loss on  $\ell(y, t)$  due to the form of  $y(\mathbf{x})$
- Can we have a smoother function such that things become easier to optimize?

 $\bullet$  An alternative: replace the  $\mathit{sign}(\cdot)$  with the sigmoid or logistic function

- An alternative: replace the  $sign(\cdot)$  with the sigmoid or logistic function
- We assumed a particular functional form: sigmoid applied to a linear function of the data

$$y(\mathbf{x}) = \sigma \left( \mathbf{w}^{\mathsf{T}} \mathbf{x} + w_0 \right)$$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

- An alternative: replace the  $sign(\cdot)$  with the sigmoid or logistic function
- We assumed a particular functional form: sigmoid applied to a linear function of the data

$$y(\mathbf{x}) = \sigma \left( \mathbf{w}^{\mathsf{T}} \mathbf{x} + w_0 \right)$$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$



- An alternative: replace the  $sign(\cdot)$  with the sigmoid or logistic function
- We assumed a particular functional form: sigmoid applied to a linear function of the data

$$y(\mathbf{x}) = \sigma \left( \mathbf{w}^{\mathsf{T}} \mathbf{x} + w_0 \right)$$

where the sigmoid is defined as

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$



• The output is a smooth function of the inputs and the weights

• We assumed a particular functional form: sigmoid applied to a linear function of the data

$$y(\mathbf{x}) = \sigma \left( \mathbf{w}^T \mathbf{x} + w_0 \right)$$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

• We assumed a particular functional form: sigmoid applied to a linear function of the data

$$y(\mathbf{x}) = \sigma \left( \mathbf{w}^{\mathsf{T}} \mathbf{x} + w_0 \right)$$

where the sigmoid is defined as

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

• One parameter per data dimension (feature)

• We assumed a particular functional form: sigmoid applied to a linear function of the data

$$y(\mathbf{x}) = \sigma \left( \mathbf{w}^T \mathbf{x} + w_0 \right)$$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

- One parameter per data dimension (feature)
- Features can be discrete or continuous

• We assumed a particular functional form: sigmoid applied to a linear function of the data

$$y(\mathbf{x}) = \sigma \left( \mathbf{w}^T \mathbf{x} + w_0 \right)$$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

- One parameter per data dimension (feature)
- Features can be discrete or continuous
- Output of the model: value  $y \in [0, 1]$

• We assumed a particular functional form: sigmoid applied to a linear function of the data

$$y(\mathbf{x}) = \sigma \left( \mathbf{w}^T \mathbf{x} + w_0 \right)$$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

- One parameter per data dimension (feature)
- Features can be discrete or continuous
- Output of the model: value  $y \in [0, 1]$
- This allows for gradient-based learning of the parameters: smoothed version of the  $\mathit{sign}(\cdot)$

## Shape of the Logistic Function

 $\bullet\,$  Let's look at how modifying w changes the function shape

# Shape of the Logistic Function

- Let's look at how modifying  ${\boldsymbol w}$  changes the function shape
- 1D example:

$$y = \sigma \left( w_1 x + w_0 \right)$$
# Shape of the Logistic Function

• Let's look at how modifying  ${\bf w}$  changes the function shape

• 1D example:

 $y = \sigma \left( w_1 x + w_0 \right)$ 



#### Show Matlab demo

• If we have a value between 0 and 1, let's use it to model the posterior

$$p(C = 0|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$
 with  $\sigma(z) = \frac{1}{1 + \exp(-z)}$ 

• If we have a value between 0 and 1, let's use it to model the posterior

$$p(C = 0|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$
 with  $\sigma(z) = \frac{1}{1 + \exp(-z)}$ 

• Substituting we have

$$p(C = 0|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^{T}\mathbf{x} - w_{0})}$$

• If we have a value between 0 and 1, let's use it to model the posterior

$$p(C = 0|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$
 with  $\sigma(z) = \frac{1}{1 + \exp(-z)}$ 

• Substituting we have

$$p(C = 0|\mathbf{x}) = \frac{1}{1 + \exp\left(-\mathbf{w}^{T}\mathbf{x} - w_{0}\right)}$$

• Supposed we have two classes, how can I compute  $p(C = 1 | \mathbf{x})$ ?

• If we have a value between 0 and 1, let's use it to model the posterior

$$p(C = 0|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$
 with  $\sigma(z) = \frac{1}{1 + \exp(-z)}$ 

• Substituting we have

$$p(C = 0|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T\mathbf{x} - w_0)}$$

- Supposed we have two classes, how can I compute  $p(C = 1 | \mathbf{x})$ ?
- Use the marginalization property of probability

$$p(C=1|\mathbf{x}) + p(C=0|\mathbf{x}) = 1$$

• If we have a value between 0 and 1, let's use it to model the posterior

$$p(C = 0|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$
 with  $\sigma(z) = \frac{1}{1 + \exp(-z)}$ 

Substituting we have

$$p(C = 0|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T\mathbf{x} - w_0)}$$

- Supposed we have two classes, how can I compute  $p(C = 1 | \mathbf{x})$ ?
- Use the marginalization property of probability

$$p(C=1|\mathbf{x})+p(C=0|\mathbf{x})=1$$

Thus (show matlab)

$$p(C = 1 | \mathbf{x}) = 1 - \frac{1}{1 + \exp\left(-\mathbf{w}^{T}\mathbf{x} - w_{0}\right)} = \frac{\exp(-\mathbf{w}^{T}\mathbf{x} - w_{0})}{1 + \exp\left(-\mathbf{w}^{T}\mathbf{x} - w_{0}\right)}$$

Assume t ∈ {0,1}, we can write the probability distribution of each of our training points p(t<sup>(1)</sup>, · · · , t<sup>(N)</sup>|x<sup>(1)</sup>, · · · x<sup>(N)</sup>)

- Assume  $t \in \{0, 1\}$ , we can write the probability distribution of each of our training points  $p(t^{(1)}, \dots, t^{(N)} | \mathbf{x}^{(1)}, \dots \mathbf{x}^{(N)})$
- Assuming that the training examples are sampled IID: independent and identically distributed

$$p(t^{(1)}, \cdots, t^{(N)} | \mathbf{x}^{(1)}, \cdots \mathbf{x}^{(N)}) = \prod_{i=1}^{N} p(t^{(i)} | \mathbf{x}^{(i)})$$

- Assume  $t \in \{0, 1\}$ , we can write the probability distribution of each of our training points  $p(t^{(1)}, \dots, t^{(N)} | \mathbf{x}^{(1)}, \dots \mathbf{x}^{(N)})$
- Assuming that the training examples are sampled IID: independent and identically distributed

$$p(t^{(1)}, \cdots, t^{(N)} | \mathbf{x}^{(1)}, \cdots \mathbf{x}^{(N)}) = \prod_{i=1}^{N} p(t^{(i)} | \mathbf{x}^{(i)})$$

• We can write each probability as

$$p(t^{(i)}|\mathbf{x}^{(i)}) = p(C = 1|\mathbf{x}^{(i)})^{t^{(i)}} p(C = 0|\mathbf{x}^{(i)})^{1-t^{(i)}}$$

- Assume  $t \in \{0, 1\}$ , we can write the probability distribution of each of our training points  $p(t^{(1)}, \dots, t^{(N)} | \mathbf{x}^{(1)}, \dots \mathbf{x}^{(N)})$
- Assuming that the training examples are sampled IID: independent and identically distributed

$$p(t^{(1)}, \cdots, t^{(N)} | \mathbf{x}^{(1)}, \cdots \mathbf{x}^{(N)}) = \prod_{i=1}^{N} p(t^{(i)} | \mathbf{x}^{(i)})$$

• We can write each probability as

$$p(t^{(i)}|\mathbf{x}^{(i)}) = p(C = 1|\mathbf{x}^{(i)})^{t^{(i)}} p(C = 0|\mathbf{x}^{(i)})^{1-t^{(i)}}$$
$$= \left(1 - p(C = 0|\mathbf{x}^{(i)})\right)^{t^{(i)}} p(C = 0|\mathbf{x}^{(i)})^{1-t^{(i)}}$$

- Assume t ∈ {0,1}, we can write the probability distribution of each of our training points p(t<sup>(1)</sup>, · · · , t<sup>(N)</sup>|x<sup>(1)</sup>, · · · x<sup>(N)</sup>)
- Assuming that the training examples are sampled IID: independent and identically distributed

$$p(t^{(1)}, \cdots, t^{(N)} | \mathbf{x}^{(1)}, \cdots \mathbf{x}^{(N)}) = \prod_{i=1}^{N} p(t^{(i)} | \mathbf{x}^{(i)})$$

• We can write each probability as

$$p(t^{(i)}|\mathbf{x}^{(i)}) = p(C = 1|\mathbf{x}^{(i)})^{t^{(i)}} p(C = 0|\mathbf{x}^{(i)})^{1-t^{(i)}}$$
$$= \left(1 - p(C = 0|\mathbf{x}^{(i)})\right)^{t^{(i)}} p(C = 0|\mathbf{x}^{(i)})^{1-t^{(i)}}$$

· We might want to learn the model, by maximizing the conditional likelihood

$$\max_{\mathbf{w}} \prod_{i=1}^{N} p(t^{(i)} | \mathbf{x}^{(i)})$$

- Assume t ∈ {0,1}, we can write the probability distribution of each of our training points p(t<sup>(1)</sup>, · · · , t<sup>(N)</sup>|x<sup>(1)</sup>, · · · x<sup>(N)</sup>)
- Assuming that the training examples are sampled IID: independent and identically distributed

$$p(t^{(1)}, \cdots, t^{(N)} | \mathbf{x}^{(1)}, \cdots \mathbf{x}^{(N)}) = \prod_{i=1}^{N} p(t^{(i)} | \mathbf{x}^{(i)})$$

• We can write each probability as

$$p(t^{(i)}|\mathbf{x}^{(i)}) = p(C = 1|\mathbf{x}^{(i)})^{t^{(i)}} p(C = 0|\mathbf{x}^{(i)})^{1-t^{(i)}}$$
$$= \left(1 - p(C = 0|\mathbf{x}^{(i)})\right)^{t^{(i)}} p(C = 0|\mathbf{x}^{(i)})^{1-t^{(i)}}$$

· We might want to learn the model, by maximizing the conditional likelihood

$$\max_{\mathbf{w}} \prod_{i=1}^{N} p(t^{(i)} | \mathbf{x}^{(i)})$$

• Convert this into a minimization so that we can write the loss function

$$p(t^{(1)}, \cdots, t^{(N)} | \mathbf{x}^{(1)}, \cdots \mathbf{x}^{(N)}) = \prod_{i=1}^{N} p(t^{(i)} | \mathbf{x}^{(i)})$$

$$p(t^{(1)}, \cdots, t^{(N)} | \mathbf{x}^{(1)}, \cdots \mathbf{x}^{(N)}) = \prod_{i=1}^{N} p(t^{(i)} | \mathbf{x}^{(i)})$$
$$= \prod_{i=1}^{N} \left( 1 - p(C = 0 | \mathbf{x}^{(i)}) \right)^{t^{(i)}} p(C = 0 | \mathbf{x}^{(i)})^{1 - t^{(i)}}$$

$$p(t^{(1)}, \cdots, t^{(N)} | \mathbf{x}^{(1)}, \cdots \mathbf{x}^{(N)}) = \prod_{i=1}^{N} p(t^{(i)} | \mathbf{x}^{(i)})$$
$$= \prod_{i=1}^{N} \left( 1 - p(C = 0 | \mathbf{x}^{(i)}) \right)^{t^{(i)}} p(C = 0 | \mathbf{x}^{(i)})^{1 - t^{(i)}}$$

• It's convenient to take the logarithm and convert the maximization into minimization by changing the sign

$$\ell_{log}(\mathbf{w}) = -\sum_{i=1}^{N} t^{(i)} \log(1 - p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})$$

$$p(t^{(1)}, \cdots, t^{(N)} | \mathbf{x}^{(1)}, \cdots \mathbf{x}^{(N)}) = \prod_{i=1}^{N} p(t^{(i)} | \mathbf{x}^{(i)})$$
$$= \prod_{i=1}^{N} \left( 1 - p(C = 0 | \mathbf{x}^{(i)}) \right)^{t^{(i)}} p(C = 0 | \mathbf{x}^{(i)})^{1 - t^{(i)}}$$

• It's convenient to take the logarithm and convert the maximization into minimization by changing the sign

$$\ell_{log}(\mathbf{w}) = -\sum_{i=1}^{N} t^{(i)} \log(1 - p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})$$

• Why is this equivalent to maximize the conditional likelihood?

$$p(t^{(1)}, \cdots, t^{(N)} | \mathbf{x}^{(1)}, \cdots \mathbf{x}^{(N)}) = \prod_{i=1}^{N} p(t^{(i)} | \mathbf{x}^{(i)})$$
$$= \prod_{i=1}^{N} \left( 1 - p(C = 0 | \mathbf{x}^{(i)}) \right)^{t^{(i)}} p(C = 0 | \mathbf{x}^{(i)})^{1 - t^{(i)}}$$

• It's convenient to take the logarithm and convert the maximization into minimization by changing the sign

$$\ell_{log}(\mathbf{w}) = -\sum_{i=1}^{N} t^{(i)} \log(1 - p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})$$

- Why is this equivalent to maximize the conditional likelihood?
- Is there a closed form solution?

$$p(t^{(1)}, \cdots, t^{(N)} | \mathbf{x}^{(1)}, \cdots \mathbf{x}^{(N)}) = \prod_{i=1}^{N} p(t^{(i)} | \mathbf{x}^{(i)})$$
$$= \prod_{i=1}^{N} \left( 1 - p(C = 0 | \mathbf{x}^{(i)}) \right)^{t^{(i)}} p(C = 0 | \mathbf{x}^{(i)})^{1 - t^{(i)}}$$

• It's convenient to take the logarithm and convert the maximization into minimization by changing the sign

$$\ell_{log}(\mathbf{w}) = -\sum_{i=1}^{N} t^{(i)} \log(1 - p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})$$

- Why is this equivalent to maximize the conditional likelihood?
- Is there a closed form solution?
- It's a convex function of w. Can we get the global optimum?

$$\min_{\mathbf{w}} \ell(\mathbf{w}) = \min_{\mathbf{w}} \left\{ -\sum_{i=1}^{N} t^{(i)} \log(1 - p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w}) \right\}$$

$$\min_{\mathbf{w}} \ell(\mathbf{w}) = \min_{\mathbf{w}} \left\{ -\sum_{i=1}^{N} t^{(i)} \log(1 - p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w}) \right\}$$

• Gradient descent: iterate and at each iteration compute steepest direction towards optimum, move in that direction, step-size  $\alpha$ 

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \alpha \frac{\partial \ell(\mathbf{w})}{\partial w_j}$$

$$\min_{\mathbf{w}} \ell(\mathbf{w}) = \min_{\mathbf{w}} \left\{ -\sum_{i=1}^{N} t^{(i)} \log(1 - \rho(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})) - \sum_{i=1}^{N} (1 - t^{(i)}) \log \rho(C = 0 | \mathbf{x}^{(i)}, \mathbf{w}) \right\}$$

• Gradient descent: iterate and at each iteration compute steepest direction towards optimum, move in that direction, step-size  $\alpha$ 

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \alpha \frac{\partial \ell(\mathbf{w})}{\partial w_j}$$

• But where is **w**?

$$\min_{\mathbf{w}} \ell(\mathbf{w}) = \min_{\mathbf{w}} \left\{ -\sum_{i=1}^{N} t^{(i)} \log(1 - \rho(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})) - \sum_{i=1}^{N} (1 - t^{(i)}) \log \rho(C = 0 | \mathbf{x}^{(i)}, \mathbf{w}) \right\}$$

• Gradient descent: iterate and at each iteration compute steepest direction towards optimum, move in that direction, step-size  $\alpha$ 

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \alpha \frac{\partial \ell(\mathbf{w})}{\partial w_j}$$

$$p(C = 0|\mathbf{x}) = \frac{1}{1 + \exp\left(-\mathbf{w}^{T}\mathbf{x} - w_{0}\right)}$$

$$\min_{\mathbf{w}} \ell(\mathbf{w}) = \min_{\mathbf{w}} \left\{ -\sum_{i=1}^{N} t^{(i)} \log(1 - p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w}) \right\}$$

• Gradient descent: iterate and at each iteration compute steepest direction towards optimum, move in that direction, step-size  $\alpha$ 

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \alpha \frac{\partial \ell(\mathbf{w})}{\partial w_j}$$

$$p(C = 0|\mathbf{x}) = \frac{1}{1 + \exp\left(-\mathbf{w}^{T}\mathbf{x} - w_{0}\right)} \qquad p(C = 1|\mathbf{x}) = \frac{\exp\left(-\mathbf{w}^{T}\mathbf{x} - w_{0}\right)}{1 + \exp\left(-\mathbf{w}^{T}\mathbf{x} - w_{0}\right)}$$

$$\min_{\mathbf{w}} \ell(\mathbf{w}) = \min_{\mathbf{w}} \left\{ -\sum_{i=1}^{N} t^{(i)} \log(1 - p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w}) \right\}$$

• Gradient descent: iterate and at each iteration compute steepest direction towards optimum, move in that direction, step-size  $\alpha$ 

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \alpha \frac{\partial \ell(\mathbf{w})}{\partial w_j}$$

$$p(C = 0|\mathbf{x}) = \frac{1}{1 + \exp\left(-\mathbf{w}^{T}\mathbf{x} - w_{0}\right)} \qquad p(C = 1|\mathbf{x}) = \frac{\exp(-\mathbf{w}^{T}\mathbf{x} - w_{0})}{1 + \exp\left(-\mathbf{w}^{T}\mathbf{x} - w_{0}\right)}$$

#### • You can write this in vector form

$$abla \ell(\mathbf{w}) = \left[\frac{\partial \ell(\mathbf{w})}{\partial w_0}, \cdots, \frac{\partial \ell(\mathbf{w})}{\partial w_k}\right]^T,$$

$$\min_{\mathbf{w}} \ell(\mathbf{w}) = \min_{\mathbf{w}} \left\{ -\sum_{i=1}^{N} t^{(i)} \log(1 - p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w}) \right\}$$

• Gradient descent: iterate and at each iteration compute steepest direction towards optimum, move in that direction, step-size  $\alpha$ 

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \alpha \frac{\partial \ell(\mathbf{w})}{\partial w_j}$$

But where is w?

$$p(C = 0|\mathbf{x}) = \frac{1}{1 + \exp\left(-\mathbf{w}^{T}\mathbf{x} - w_{0}\right)} \qquad p(C = 1|\mathbf{x}) = \frac{\exp(-\mathbf{w}^{T}\mathbf{x} - w_{0})}{1 + \exp\left(-\mathbf{w}^{T}\mathbf{x} - w_{0}\right)}$$

#### • You can write this in vector form

$$\nabla \ell(\mathbf{w}) = \left[\frac{\partial \ell(\mathbf{w})}{\partial w_0}, \cdots, \frac{\partial \ell(\mathbf{w})}{\partial w_k}\right]^T, \quad \text{and} \quad \triangle(\mathbf{w}) = -\alpha \bigtriangledown \ell(\mathbf{w})$$

• The log likelihood is

$$\ell_{log-loss}(\mathbf{w}) = -\sum_{i=1}^{N} t^{(i)} \log p(C = 1 | \mathbf{x}^{(i)}, \mathbf{w}) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})$$

• The log likelihood is

$$\ell_{log-loss}(\mathbf{w}) = -\sum_{i=1}^{N} t^{(i)} \log p(C = 1 | \mathbf{x}^{(i)}, \mathbf{w}) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})$$

where the probabilities are

$$p(C = 0 | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-z)} \qquad p(C = 1 | \mathbf{x}, \mathbf{w}) = \frac{\exp(-z)}{1 + \exp(-z)}$$

• The log likelihood is

$$\ell_{log-loss}(\mathbf{w}) = -\sum_{i=1}^{N} t^{(i)} \log p(C = 1 | \mathbf{x}^{(i)}, \mathbf{w}) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})$$

where the probabilities are

$$p(C=0|\mathbf{x},\mathbf{w}) = \frac{1}{1+\exp(-z)} \qquad p(C=1|\mathbf{x},\mathbf{w}) = \frac{\exp(-z)}{1+\exp(-z)}$$

and  $z = \mathbf{w}^T \mathbf{x} + w_0$ 

• The log likelihood is

$$\ell_{log-loss}(\mathbf{w}) = -\sum_{i=1}^{N} t^{(i)} \log p(C = 1 | \mathbf{x}^{(i)}, \mathbf{w}) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})$$

where the probabilities are

$$p(C = 0 | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-z)} \qquad p(C = 1 | \mathbf{x}, \mathbf{w}) = \frac{\exp(-z)}{1 + \exp(-z)}$$

and  $z = \mathbf{w}^T \mathbf{x} + w_0$ 

• We can simplify

$$\ell(\mathbf{w}) = \sum_{i} t^{(i)} \log(1 + \exp(-z^{(i)})) + \sum_{i} t^{(i)} z^{(i)} + \sum_{i} (1 - t^{(i)}) \log(1 + \exp(-z^{(i)}))$$

• The log likelihood is

$$\ell_{log-loss}(\mathbf{w}) = -\sum_{i=1}^{N} t^{(i)} \log p(C = 1 | \mathbf{x}^{(i)}, \mathbf{w}) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})$$

where the probabilities are

$$p(C = 0 | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-z)} \qquad p(C = 1 | \mathbf{x}, \mathbf{w}) = \frac{\exp(-z)}{1 + \exp(-z)}$$

and  $z = \mathbf{w}^T \mathbf{x} + w_0$ 

• We can simplify

$$\begin{split} \ell(\mathbf{w}) &= \sum_{i} t^{(i)} \log(1 + \exp(-z^{(i)})) + \sum_{i} t^{(i)} z^{(i)} + \sum_{i} (1 - t^{(i)}) \log(1 + \exp(-z^{(i)})) \\ &= \sum_{i} \log(1 + \exp(-z^{(i)})) + \sum_{i} t^{(i)} z^{(i)} \end{split}$$

• The log likelihood is

$$\ell_{log-loss}(\mathbf{w}) = -\sum_{i=1}^{N} t^{(i)} \log p(C = 1 | \mathbf{x}^{(i)}, \mathbf{w}) - \sum_{i=1}^{N} (1 - t^{(i)}) \log p(C = 0 | \mathbf{x}^{(i)}, \mathbf{w})$$

where the probabilities are

$$p(C = 0 | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-z)} \qquad p(C = 1 | \mathbf{x}, \mathbf{w}) = \frac{\exp(-z)}{1 + \exp(-z)}$$

and  $z = \mathbf{w}^T \mathbf{x} + w_0$ 

• We can simplify

$$\ell(\mathbf{w}) = \sum_{i} t^{(i)} \log(1 + \exp(-z^{(i)})) + \sum_{i} t^{(i)} z^{(i)} + \sum_{i} (1 - t^{(i)}) \log(1 + \exp(-z^{(i)}))$$
  
= 
$$\sum_{i} \log(1 + \exp(-z^{(i)})) + \sum_{i} t^{(i)} z^{(i)}$$

• Now it's easy to take derivatives

$$\ell(\mathbf{w}) = \sum_{i} t^{(i)} z^{(i)} + \sum_{i} \log(1 + \exp(-z^{(i)}))$$

• Now it's easy to take derivatives

$$\ell(\mathbf{w}) = \sum_{i} t^{(i)} z^{(i)} + \sum_{i} \log(1 + \exp(-z^{(i)}))$$

- Now it's easy to take derivatives
- Remember  $z = \mathbf{w}^T \mathbf{x} + w_0$

$$\ell(\mathbf{w}) = \sum_{i} t^{(i)} z^{(i)} + \sum_{i} \log(1 + \exp(-z^{(i)}))$$

- Now it's easy to take derivatives
- Remember  $z = \mathbf{w}^T \mathbf{x} + w_0$

$$\frac{\partial \ell(\mathbf{w})}{\partial w_j} = \sum_i t^{(i)} x_j^{(i)} - x_j^{(i)} \cdot \frac{\exp(-z^{(i)})}{1 + \exp(-z^{(i)})}$$

$$\ell(\mathbf{w}) = \sum_{i} t^{(i)} z^{(i)} + \sum_{i} \log(1 + \exp(-z^{(i)}))$$

- Now it's easy to take derivatives
- Remember  $z = \mathbf{w}^T \mathbf{x} + w_0$

$$\frac{\partial \ell(\mathbf{w})}{\partial w_j} = \sum_i t^{(i)} x_j^{(i)} - x_j^{(i)} \cdot \frac{\exp(-z^{(i)})}{1 + \exp(-z^{(i)})}$$

• What's 
$$x_j^{(i)}$$
?

$$\ell(\mathbf{w}) = \sum_{i} t^{(i)} z^{(i)} + \sum_{i} \log(1 + \exp(-z^{(i)}))$$

- Now it's easy to take derivatives
- Remember  $z = \mathbf{w}^T \mathbf{x} + w_0$

$$\frac{\partial \ell(\mathbf{w})}{\partial w_j} = \sum_i t^{(i)} x_j^{(i)} - x_j^{(i)} \cdot \frac{\exp(-z^{(i)})}{1 + \exp(-z^{(i)})}$$

- What's  $x_j^{(i)}$ ?
- And simplifying

$$\frac{\partial \ell(\mathbf{w})}{\partial w_j} = \sum_i x_j^{(i)} \left( t^{(i)} - p(C = 1 | \mathbf{x}^{(i)}) \right)$$
## Updates

$$\ell(\mathbf{w}) = \sum_{i} t^{(i)} z^{(i)} + \sum_{i} \log(1 + \exp(-z^{(i)}))$$

- Now it's easy to take derivatives
- Remember  $z = \mathbf{w}^T \mathbf{x} + w_0$

$$\frac{\partial \ell(\mathbf{w})}{\partial w_j} = \sum_i t^{(i)} x_j^{(i)} - x_j^{(i)} \cdot \frac{\exp(-z^{(i)})}{1 + \exp(-z^{(i)})}$$

- What's  $x_j^{(i)}$ ?
- And simplifying

$$\frac{\partial \ell(\mathbf{w})}{\partial w_j} = \sum_i x_j^{(i)} \left( t^{(i)} - p(C = 1 | \mathbf{x}^{(i)}) \right)$$

• Don't get confused with indexes: *j* for the weight that we are updating and *i* for the training example

## Updates

$$\ell(\mathbf{w}) = \sum_{i} t^{(i)} z^{(i)} + \sum_{i} \log(1 + \exp(-z^{(i)}))$$

- Now it's easy to take derivatives
- Remember  $z = \mathbf{w}^T \mathbf{x} + w_0$

$$\frac{\partial \ell(\mathbf{w})}{\partial w_j} = \sum_i t^{(i)} x_j^{(i)} - x_j^{(i)} \cdot \frac{\exp(-z^{(i)})}{1 + \exp(-z^{(i)})}$$

- What's  $x_j^{(i)}$ ?
- And simplifying

$$\frac{\partial \ell(\mathbf{w})}{\partial w_j} = \sum_i x_j^{(i)} \left( t^{(i)} - p(C = 1 | \mathbf{x}^{(i)}) \right)$$

- Don't get confused with indexes: *j* for the weight that we are updating and *i* for the training example
- Logistic regression has linear decision boundary

Urtasun & Zemel (UofT)

### Logistic regression vs least squares





If the right answer is 1 and the model says 1.5, it loses, so it changes the boundary to avoid being "too correct" (tilts aways from outliers) • We can also look at

 $p(\mathbf{w}|\{t\},\{\mathbf{x}\}) \propto p(\{t\}|\{\mathbf{x}\},\mathbf{w}) p(\mathbf{w})$ with  $\{t\} = (t^1, \cdots, t^N)$ , and  $\{\mathbf{x}\} = (\mathbf{x}^1, \cdots, \mathbf{x}^N)$  • We can also look at

$$p(\mathbf{w}|\{t\}, \{\mathbf{x}\}) \propto p(\{t\}|\{\mathbf{x}\}, \mathbf{w}) p(\mathbf{w})$$
with  $\{t\} = (t^1, \cdots, t^N)$ , and  $\{\mathbf{x}\} = (\mathbf{x}^1, \cdots, \mathbf{x}^N)$ 

 $\bullet\,$  We can define priors on parameters  ${\bf w}$ 

• We can also look at

$$p(\mathbf{w}|\{t\}, \{\mathbf{x}\}) \propto p(\{t\}|\{\mathbf{x}\}, \mathbf{w}) p(\mathbf{w})$$
with  $\{t\} = (t^1, \cdots, t^N)$ , and  $\{\mathbf{x}\} = (\mathbf{x}^1, \cdots, \mathbf{x}^N)$ 

- ${\ensuremath{\bullet}}$  We can define priors on parameters  ${\ensuremath{w}}$
- This is a form of regularization

We can also look at

 $p(\mathbf{w}|\{t\}, \{\mathbf{x}\}) \propto p(\{t\}|\{\mathbf{x}\}, \mathbf{w}) p(\mathbf{w})$ 

with  $\{t\} = (t^1, \cdots, t^N)$ , and  $\{\mathbf{x}\} = (\mathbf{x}^1, \cdots, \mathbf{x}^N)$ 

- We can define priors on parameters w
- This is a form of regularization
- Helps avoid large weights and over fitting

$$\max_{\mathbf{w}} \log \left[ p(\mathbf{w}) \prod_{i} p(t^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \right]$$

We can also look at

 $p(\mathbf{w}|\{t\}, \{\mathbf{x}\}) \propto p(\{t\}|\{\mathbf{x}\}, \mathbf{w}) p(\mathbf{w})$ 

with  $\{t\} = (t^1, \cdots, t^N)$ , and  $\{\mathbf{x}\} = (\mathbf{x}^1, \cdots, \mathbf{x}^N)$ 

- We can define priors on parameters w
- This is a form of regularization
- Helps avoid large weights and over fitting

$$\max_{\mathbf{w}} \log \left[ p(\mathbf{w}) \prod_{i} p(t^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \right]$$

• What's  $p(\mathbf{w})$ ?

• For example, define prior: normal distribution, zero mean and identity covariance  $p(\mathbf{w}) = \mathcal{N}(0, \lambda \mathbf{I})$ 

- For example, define prior: normal distribution, zero mean and identity covariance  $p(\mathbf{w}) = \mathcal{N}(0, \lambda \mathbf{I})$
- This prior pushes parameters towards zero

- For example, define prior: normal distribution, zero mean and identity covariance  $p(\mathbf{w}) = \mathcal{N}(0, \lambda \mathbf{I})$
- This prior pushes parameters towards zero
- Including this prior the new gradient is

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \alpha \frac{\partial \ell(\mathbf{w})}{\partial w_j} - \alpha \lambda w_j^{(t)}$$

- For example, define prior: normal distribution, zero mean and identity covariance  $p(\mathbf{w}) = \mathcal{N}(0, \lambda \mathbf{I})$
- This prior pushes parameters towards zero
- Including this prior the new gradient is

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \alpha \frac{\partial \ell(\mathbf{w})}{\partial w_j} - \alpha \lambda w_j^{(t)}$$

• How do we decide the best value of  $\lambda$ ?

• We can divide the set of training examples into two disjoint sets: training and validation

- We can divide the set of training examples into two disjoint sets: training and validation
- Use the first set (i.e., training) to estimate the weights  ${\bf w}$  for different values of  $\lambda$

- We can divide the set of training examples into two disjoint sets: training and validation
- Use the first set (i.e., training) to estimate the weights  ${\bf w}$  for different values of  $\lambda$
- Use the second set (i.e., validation) to estimate the best λ, by evaluating how well the classifier does in this second set

- We can divide the set of training examples into two disjoint sets: training and validation
- Use the first set (i.e., training) to estimate the weights  ${\bf w}$  for different values of  $\lambda$
- Use the second set (i.e., validation) to estimate the best λ, by evaluating how well the classifier does in this second set
- This test how well you generalized to unseen data

- We can divide the set of training examples into two disjoint sets: training and validation
- Use the first set (i.e., training) to estimate the weights  ${\bf w}$  for different values of  $\lambda$
- Use the second set (i.e., validation) to estimate the best λ, by evaluating how well the classifier does in this second set
- This test how well you generalized to unseen data
- The parameter  $\lambda$  is the importance of the regularization, and it's a  $\mbox{hyper}$  parameter

• We use *p* observations as the validation set and the remaining observations as the training set.

- We use *p* observations as the validation set and the remaining observations as the training set.
- This is repeated on all ways to cut the original training set.

- We use *p* observations as the validation set and the remaining observations as the training set.
- This is repeated on all ways to cut the original training set.
- It requires  $C_n^p$  for a set of n examples

- We use *p* observations as the validation set and the remaining observations as the training set.
- This is repeated on all ways to cut the original training set.
- It requires  $C_n^p$  for a set of *n* examples

• Leave-1-out cross-validation: When p = 1, does not have this problem

- We use *p* observations as the validation set and the remaining observations as the training set.
- This is repeated on all ways to cut the original training set.
- It requires  $C_n^p$  for a set of *n* examples
- Leave-1-out cross-validation: When p = 1, does not have this problem
- k-fold cross-validation:

- We use *p* observations as the validation set and the remaining observations as the training set.
- This is repeated on all ways to cut the original training set.
- It requires  $C_n^p$  for a set of n examples
- Leave-1-out cross-validation: When p = 1, does not have this problem
- k-fold cross-validation:
  - The training set is randomly partitioned into k equal size subsamples.

- We use *p* observations as the validation set and the remaining observations as the training set.
- This is repeated on all ways to cut the original training set.
- It requires  $C_n^p$  for a set of n examples
- Leave-1-out cross-validation: When p = 1, does not have this problem
- k-fold cross-validation:
  - The training set is randomly partitioned into k equal size subsamples.
  - Of the k subsamples, a single subsample is retained as the validation data for testing the model, and the remaining k 1 subsamples are used as training data.

- We use *p* observations as the validation set and the remaining observations as the training set.
- This is repeated on all ways to cut the original training set.
- It requires  $C_n^p$  for a set of n examples
- Leave-1-out cross-validation: When p = 1, does not have this problem

#### • k-fold cross-validation:

- The training set is randomly partitioned into k equal size subsamples.
- Of the k subsamples, a single subsample is retained as the validation data for testing the model, and the remaining k 1 subsamples are used as training data.
- The cross-validation process is then repeated k times (the folds).

- We use *p* observations as the validation set and the remaining observations as the training set.
- This is repeated on all ways to cut the original training set.
- It requires  $C_n^p$  for a set of n examples
- Leave-1-out cross-validation: When p = 1, does not have this problem

#### • k-fold cross-validation:

- The training set is randomly partitioned into k equal size subsamples.
- Of the k subsamples, a single subsample is retained as the validation data for testing the model, and the remaining k 1 subsamples are used as training data.
- The cross-validation process is then repeated *k* times (the folds).
- The k results from the folds can then be averaged (or otherwise combined) to produce a single estimation

### Metrics

How to evaluate how good my classifier is?

• Precision: is the fraction of retrieved instances that are relevant

$$P = \frac{TP}{TP + FN} = \frac{TP}{P}$$

### Metrics

How to evaluate how good my classifier is?

• Precision: is the fraction of retrieved instances that are relevant

$$P = \frac{TP}{TP + FN} = \frac{TP}{P}$$

• Recall: is the fraction of relevant instances that are retrieved

$$R = \frac{TP}{TP + FP}$$

### Metrics

How to evaluate how good my classifier is?

• Precision: is the fraction of retrieved instances that are relevant

$$P = \frac{TP}{TP + FN} = \frac{TP}{P}$$

• Recall: is the fraction of relevant instances that are retrieved

$$R = \frac{TP}{TP + FP}$$

• F1 score: harmonic mean of precision and recall

$$F1 = 2\frac{P \cdot R}{P + R}$$

# More on Metrics

How to evaluate how good my classifier is?

- Precision: is the fraction of retrieved instances that are relevant
- Recall: is the fraction of relevant instances that are retrieved

## More on Metrics

How to evaluate how good my classifier is?

- Precision: is the fraction of retrieved instances that are relevant
- Recall: is the fraction of relevant instances that are retrieved
- Precision Recall Curve



## More on Metrics

How to evaluate how good my classifier is?

- Precision: is the fraction of retrieved instances that are relevant
- Recall: is the fraction of relevant instances that are retrieved
- Precision Recall Curve



• Average Precision (AP): mean under the curve

# Classification: Oranges and Lemons



## Classification: Oranges and Lemons



• Linear classification means that the part that adapts is linear (just like linear regression)

$$z(x) = \mathbf{w}^T \mathbf{x} + w_0$$

with adaptive **w**, w<sub>0</sub>

• The adaptive part is follow by a non-linearity

$$y(\mathbf{x}) = f(z(\mathbf{x}))$$

• What *f* have we seen?
### Classification as Induction



#### • Alternative to parametric model is non-parametric

- Alternative to parametric model is non-parametric
- Simple methods for approximating discrete-valued or real-valued target functions (classification or regression problems)

- Alternative to parametric model is non-parametric
- Simple methods for approximating discrete-valued or real-valued target functions (classification or regression problems)
- Learning amounts to simply storing training data

- Alternative to parametric model is non-parametric
- Simple methods for approximating discrete-valued or real-valued target functions (classification or regression problems)
- Learning amounts to simply storing training data
- Test instances classified using similar training instances

- Alternative to parametric model is non-parametric
- Simple methods for approximating discrete-valued or real-valued target functions (classification or regression problems)
- Learning amounts to simply storing training data
- Test instances classified using similar training instances
- Embodies often sensible underlying assumptions:
  - Output varies smoothly with input
  - Data occupies sub-space of high-dimensional input space

• Assume training examples correspond to points in d-dimensional Euclidean space

- Assume training examples correspond to points in d-dimensional Euclidean space
- Target function value for new query estimated from known value of nearest training example(s)

- Assume training examples correspond to points in d-dimensional Euclidean space
- Target function value for new query estimated from known value of nearest training example(s)
- Distance typically defined to be Euclidean:

$$||\mathbf{x}^{(a)} - \mathbf{x}^{(b)}|| = \sqrt{\sum_{j=1}^{d} (x_j^{(a)} - (x_j^{(b)})^2)^2}$$

- Assume training examples correspond to points in d-dimensional Euclidean space
- Target function value for new query estimated from known value of nearest training example(s)
- Distance typically defined to be Euclidean:

$$||\mathbf{x}^{(a)} - \mathbf{x}^{(b)}|| = \sqrt{\sum_{j=1}^{d} (x_j^{(a)} - (x_j^{(b)})^2)^2}$$

Algorithm

General Ansatz (a) find example (x\*, t\*) closest to the test instance x<sup>(q)</sup>
 Output y<sup>(q)</sup> = t\*

- Assume training examples correspond to points in d-dimensional Euclidean space
- Target function value for new query estimated from known value of nearest training example(s)
- Distance typically defined to be Euclidean:

$$||\mathbf{x}^{(a)} - \mathbf{x}^{(b)}|| = \sqrt{\sum_{j=1}^{d} (x_j^{(a)} - (x_j^{(b)})^2)^2}$$

Algorithm

- find example (x\*, t\*) closest to the test instance x<sup>(q)</sup>
  output y<sup>(q)</sup> = t\*
- Note: we don't need to compute the square root. Why?

### Nearest Neighbors Decision Boundaries

- Nearest neighbor algorithm does not explicitly compute decision boundaries, but these can be inferred
- Decision boundaries: Voronoi diagram visualization
  - show how input space divided into classes
  - each line segment is equidistant between two points of opposite classes



 $\bullet$  Nearest neighbors sensitive to mis-labeled data ("class noise" )  $\to$  smooth by having k nearest neighbors vote

- $\bullet$  Nearest neighbors sensitive to mis-labeled data ("class noise")  $\to$  smooth by having k nearest neighbors vote
- Algorithm:
  - **1** find k examples  $\{\mathbf{x}^{(i)}, t^{(i)}\}$  closest to the test instance  $\mathbf{x}$
  - elassification output is majority class

$$y = \arg \max_{t^{(z)}} \sum_{r=1}^k \delta(t^{(z)}, t^{(r)})$$

 $\bullet\,$  Some attributes have larger ranges, so are treated as more important  $\to\,$  normalize scale

- $\bullet\,$  Some attributes have larger ranges, so are treated as more important  $\to\,$  normalize scale
- Irrelevant, correlated attributes add noise to distance measure  $\rightarrow$  eliminate some attributes, or vary and possibly adapt weight of attributes

- $\bullet\,$  Some attributes have larger ranges, so are treated as more important  $\to\,$  normalize scale
- Irrelevant, correlated attributes add noise to distance measure  $\rightarrow$  eliminate some attributes, or vary and possibly adapt weight of attributes
- Non-metric attributes (symbols)  $\rightarrow$  Hamming distance

- $\bullet\,$  Some attributes have larger ranges, so are treated as more important  $\to\,$  normalize scale
- Irrelevant, correlated attributes add noise to distance measure  $\rightarrow$  eliminate some attributes, or vary and possibly adapt weight of attributes
- Non-metric attributes (symbols)  $\rightarrow$  Hamming distance
- Brute-force approach: calculate Euclidean distance to test point from each stored point, keep closest:  $O(dn^2) \rightarrow$  reduce computational burden:
  - Use subset of dimensions
  - ② Use subset of examples

- $\bullet\,$  Some attributes have larger ranges, so are treated as more important  $\to\,$  normalize scale
- Irrelevant, correlated attributes add noise to distance measure  $\rightarrow$  eliminate some attributes, or vary and possibly adapt weight of attributes
- Non-metric attributes (symbols)  $\rightarrow$  Hamming distance
- Brute-force approach: calculate Euclidean distance to test point from each stored point, keep closest:  $O(dn^2) \rightarrow$  reduce computational burden:
  - Use subset of dimensions
  - ② Use subset of examples
    - Remove examples that lie within Voronoi region

- $\bullet\,$  Some attributes have larger ranges, so are treated as more important  $\to\,$  normalize scale
- Irrelevant, correlated attributes add noise to distance measure  $\rightarrow$  eliminate some attributes, or vary and possibly adapt weight of attributes
- Non-metric attributes (symbols)  $\rightarrow$  Hamming distance
- Brute-force approach: calculate Euclidean distance to test point from each stored point, keep closest:  $O(dn^2) \rightarrow$  reduce computational burden:
  - Use subset of dimensions
  - ② Use subset of examples
    - Remove examples that lie within Voronoi region
    - Form efficient search tree (kd-tree), use Hashing (LSH), etc

### Decision Boundary K-NN





• Single parameter (k)  $\rightarrow$  how do we set it?



- Single parameter  $(k) \rightarrow how do we set it?$
- Naturally forms complex decision boundaries; adapts to data density



- Single parameter  $(k) \rightarrow how do we set it?$
- Naturally forms complex decision boundaries; adapts to data density
- Problems:
  - Sensitive to class noise.
  - Sensitive to dimensional scales.
  - Distances are less meaningful in high dimensions
  - Scales with number of examples



- Single parameter (k)  $\rightarrow$  how do we set it?
- Naturally forms complex decision boundaries; adapts to data density
- Problems:
  - Sensitive to class noise.
  - Sensitive to dimensional scales.
  - Distances are less meaningful in high dimensions
  - Scales with number of examples
- Inductive Bias: What kind of decision boundaries do we expect to find?