# STA 314: Statistical Methods for Machine Learning I Lecture 6 - Linear Classification II

Chris J. Maddison

University of Toronto

## **Today**

- Support vector machines, elegant binary linear classifiers that generalize very well.
- Multiclass classification: predicting a discrete(> 2)-valued target.
- Stochastic gradient descent, which lets us scale up gradient descent to large data sets.

Intro ML(UofT) STA314-Lec6 2/44

## Binary Classification with a Linear Model (Small Change)

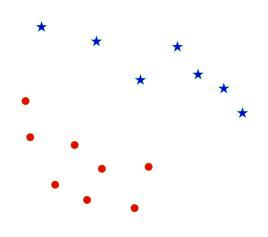
- Binary classification: predicting a target with two values
- Targets (small change from last week):  $t \in \{-1, +1\}$
- Linear model (small change from last week):

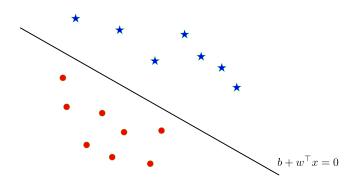
$$z = \mathbf{w}^{\top} \mathbf{x} + b$$
$$y = sign(z)$$

- This is an equivalent formulation of binary linear classification.
- Last week we considered how to get any w and b that minimized the cost on the training set.
- $\bullet$  Question: How should we choose  $\mathbf{w}$  and b to get the best generalization?

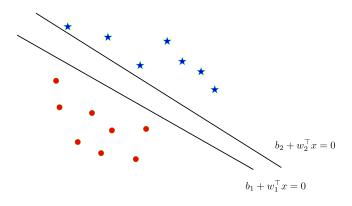
Intro ML(UofT) STA314-Lec6 3 / 44

Suppose we are given these data points from two different classes and want to find a linear classifier that separates them.



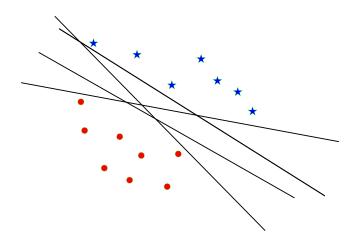


- The decision boundary looks like a line because  $\mathbf{x} \in \mathbb{R}^2$ , but think about it as a D-1 dimensional hyperplane.
- Recall that a hyperplane is described by points  $\mathbf{x} \in \mathbb{R}^D$  such that  $f(\mathbf{x}) = \mathbf{w}^T x + b = 0$ .



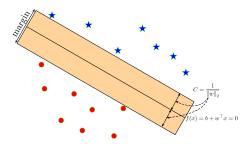
• There are multiple separating hyperplanes, described by different parameters  $(\mathbf{w}, b)$ .

Intro ML(UofT) STA314-Lec6 6 / 44



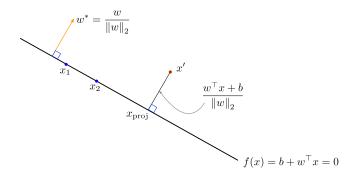
## Optimal Separating Hyperplane

Optimal Separating Hyperplane: A hyperplane that separates two classes and maximizes the distance to the closest point from either class, i.e., maximize the margin of the classifier.



Intuitively, ensuring that a classifier is not too close to any data points leads to better generalization on the test data.

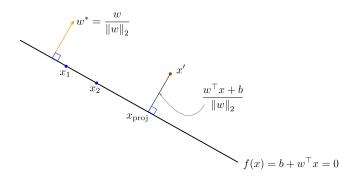
## Geometry of Points and Planes



• Recall that the decision hyperplane is orthogonal (perpendicular) to  $\mathbf{w}$ . I.e., for any two points  $\mathbf{x}_1$  and  $\mathbf{x}_2$  on the decision hyperplane we have that  $\mathbf{w}^{\top}(\mathbf{x}_1 - \mathbf{x}_2) = 0$ .

Intro ML(UofT) STA314-Lec6 9 / 44

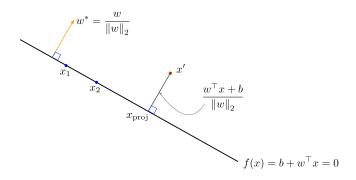
## Geometry of Points and Planes



- The vector  $\mathbf{w}^* = \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$  is a unit vector pointing in the same direction as  $\mathbf{w}$ .
- The same hyperplane could equivalently be defined in terms of  $\mathbf{w}^*$ .

Intro ML (UofT) STA314-Lec6 10 / 44

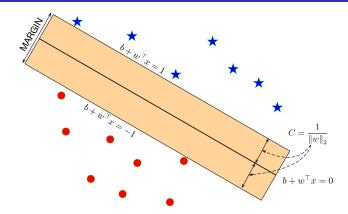
## Geometry of Points and Planes



• To get the distance from a point  $\mathbf{x}$  to the hyperplane, take the closest point  $\mathbf{x}_{\mathrm{proj}}$  on the hyperplane and project  $\mathbf{x} - \mathbf{x}_{\mathrm{proj}}$  onto  $\mathbf{w} / \|\mathbf{w}\|_2$ :

$$\left| \left( \mathbf{x} - \mathbf{x}_{\text{proj}} \right)^{\mathsf{T}} \frac{\mathbf{w}}{\left\| \mathbf{w} \right\|_{2}} \right| = \frac{\left| \mathbf{x}^{\mathsf{T}} \mathbf{w} - \mathbf{x}_{\text{proj}}^{\mathsf{T}} \mathbf{w} \right|}{\left\| \mathbf{w} \right\|_{2}} = \frac{\left| \mathbf{x}^{\mathsf{T}} \mathbf{w} + b \right|}{\left\| \mathbf{w} \right\|_{2}}$$

Intro ML (UofT) STA314-Lec6 11/44



• Now consider the two parallel hyperplanes

$$\mathbf{w}^{\mathsf{T}}\mathbf{x} + b = 1$$
  $\mathbf{w}^{\mathsf{T}}\mathbf{x} + b = -1$ 

• Using the distance formula, can see that the margin is  $2/\|\mathbf{w}\|_2$ .

Intro ML (UofT) STA314-Lec6 12 / 44

Recall: to correctly classify all points we require that

$$sign(\mathbf{w}^{\top}\mathbf{x}^{(i)} + b) = t^{(i)}$$
 for all  $i$ 

• Let's impose a stronger requirement: correctly classify all points *and* prevent them from falling in the margin.

$$\mathbf{w}^{\mathsf{T}} \mathbf{x}^{(i)} + b \ge 1$$
 if  $t^{(i)} = 1$   
 $\mathbf{w}^{\mathsf{T}} \mathbf{x}^{(i)} + b \le -1$  if  $t^{(i)} = -1$ 

• This is equivalent to

$$t^{(i)} \left( \mathbf{w}^{\top} \mathbf{x}^{(i)} + b \right) \ge 1$$
 for all  $i$ 

which we call the margin constraints.

- Now, we want to pick w, b that maximize the size of the margin (the region where we do not allow points to fall), while ensuring all points are correctly classified.
  - ▶ Margin has width  $2/||\mathbf{w}||_2$ , so maximizing this is equivalent to minimizing  $||\mathbf{w}||_2^2$ .
- This leads to the max-margin objective:

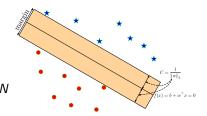
$$\min_{\mathbf{w},b} \|\mathbf{w}\|_{2}^{2}$$
s.t.  $t^{(i)}(\mathbf{w}^{\top}\mathbf{x}^{(i)} + b) \ge 1$   $i = 1,..., N$ 

Intro ML(UofT) STA314-Lec6 14 / 44

Max-margin objective:

$$\min_{\mathbf{w},b} \|\mathbf{w}\|_2^2$$

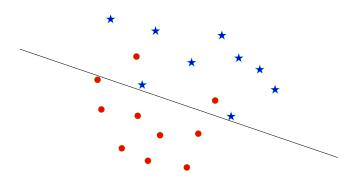
s.t. 
$$t^{(i)}(\mathbf{w}^{\top}\mathbf{x}^{(i)} + b) \ge 1$$
  $i = 1, ..., N$ 



- Observe: if the margin constraint is not tight for  $\mathbf{x}^{(i)}$ , we could remove it from the training set and the optimal  $\mathbf{w}$  would be the same.
- The important training examples are the ones with algebraic margin 1, and are called support vectors.
- Hence, this algorithm is called the (hard) Support Vector Machine (SVM) (or Support Vector Classifier).
- SVM-like algorithms are often called max-margin or large-margin.

## Non-Separable Data Points

How can we apply the max-margin principle if the data are **not** linearly separable?



## Hinge Loss

- There is an elegant relaxation of the max-margin objective called soft-margin SVM for the non-separable case. We won't cover it carefully, but let's motivate it.
- We can measure the extent to which  $\mathbf{x}^{(i)}$  violates its margin constraint by the magnitude of

$$1 - t^{(i)}(\mathbf{w}^{\top}\mathbf{x}^{(i)} + b)$$

- ▶ If this is very positive, then  $t^{(i)}(\mathbf{w}^{\top}\mathbf{x}^{(i)} + b) < 1$ .
- ▶ If this is not positive, then  $t^{(i)}(\mathbf{w}^{\top}\mathbf{x}^{(i)} + b) \ge 1$ .

## Hinge Loss

- We want 1 t<sup>(i)</sup>(w<sup>T</sup>x<sup>(i)</sup> + b) to be small, but if it is negative, we don't care how negative it is.
- This motivates the hinge loss for  $z^{(i)} = \mathbf{w}^{\mathsf{T}} \mathbf{x}^{(i)} + b$ :

$$L_{\text{hinge}}(z^{(i)}) = \max(0, 1 - t^{(i)}z^{(i)})$$

## Soft-margin SVM

 Soft-margin SVM minimizes the average hinge losses plus the norm of the weights, where  $z^{(i)} = \mathbf{w}^{\top} \mathbf{x}^{(i)} + b$ :

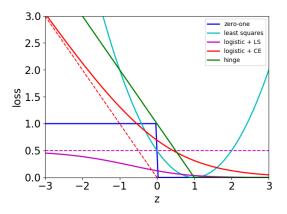
$$\min_{\mathbf{w},b} \sum_{i=1}^{N} \frac{1}{N} \max(0, 1 - t^{(i)} z^{(i)}) + \lambda \|\mathbf{w}\|_{2}^{2}$$

 Hence, the soft-margin SVM can be seen as a linear classifier with hinge loss and an  $L_2$  regularizer.

STA314-Lec6 19 / 44

## Revisiting Loss Functions for Classification

Hinge loss compared with other loss functions



#### SVMs: What we Left Out

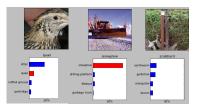
#### What we left out:

- How to fit w for the max-margin SVM:
  - ▶ One option: gradient descent
- Classic results from learning theory show that a large margin implies good generalization.

#### Multiclass Classification

- So far, we've only talked about binary linear classification.
- Classification tasks with more than two categories:





#### Multiclass Classification

- Targets form a discrete set  $\{1, \ldots, K\}$ .
- It's often more convenient to represent them as one-hot vectors, or a one-of-K encoding:

$$\mathbf{t} = (0, \dots, 0, 1, 0, \dots, 0) \in \mathbb{R}^{K}$$
entry  $k$  is 1

 We will build multiclass linear classifiers by generalizing binary linear classifiers and logistic regression (not SVMs).

### Multiclass Linear Classification

- We can start with a linear function of the inputs.
- D input dimensions and K output dimensions, so we need  $K \times D$  weights, which we arrange as a weight matrix  $\mathbf{W} \in \mathbb{R}^{K \times D}$ .
- Also, we have a vector  $\mathbf{b} \in \mathbb{R}^K$  of bias parameters.
- A linear function of an input  $\mathbf{x} \in \mathbb{R}^D$ :

$$z_k = \sum_{j=1}^{D} w_{kj} x_j + b_k$$
 for  $k = 1, 2, ..., K$ 

- Eliminate the bias parameters by taking  $\mathbf{W} \in \mathbb{R}^{K \times (D+1)}$  and adding a dummy variable  $x_0 = 1$ .
- So, vectorized we have the vector  $\mathbf{z} \in \mathbb{R}^K$ :

$$z = Wx + b$$
 or with dummy  $x_0 = 1$   $z = Wx$ 

### Multiclass Linear Classification

- How can we turn this linear prediction into a one-hot prediction?
- We can interpret the magnitude of  $z_k$  as an measure of how much the model prefers k as its prediction.
- If we do this, we should set

$$y_i = \begin{cases} 1 & i = \arg\max_{k=1}^K z_k \\ 0 & \text{otherwise} \end{cases}$$

• Exercise: how does the case of K = 2 relate to the prediction rule in binary linear classifiers?

## Softmax Regression

- As with binary classification, we need to soften our predictions for the sake of optimization.
- We want predictions that are like probabilities, i.e.,  $0 \le y_k \le 1$  and  $\sum_k y_k = 1$ .
- We can use the softmax function, a multivariable generalization of the logistic function:

$$y_k = \operatorname{softmax}(\mathbf{z})_k = \frac{e^{z_k}}{\sum_{k'} e^{z_{k'}}}$$

- Outputs can be interpreted as probabilities (positive and sum to 1)
- ▶ If  $z_k$  is much larger than the others, then  $\operatorname{softmax}(\mathbf{z})_k \approx 1$  and it behaves like argmax.
- **Exercise:** how does the case of K = 2 relate to the logistic function?
- The inputs  $z_k$  are called the logits.

## Softmax Regression

• If a model outputs a vector  $\mathbf{y} \in \mathbb{R}^K$  of class probabilities, we can use cross-entropy as the loss function:

$$L_{\text{CE}}(\mathbf{y}, \mathbf{t}) = -\sum_{k=1}^{K} t_k \log y_k$$
$$= -\mathbf{t}^{\top} (\log \mathbf{y}),$$

where the log is applied elementwise.

 Just like with logistic regression, we typically combine the softmax and cross-entropy into a softmax-cross-entropy function.

## Softmax Regression

• Softmax regression (with dummy  $x_0 = 1$ ):

$$\begin{aligned} \mathbf{z} &= \mathbf{W} \mathbf{x} \\ \mathbf{y} &= \operatorname{softmax}(\mathbf{z}) \\ L_{\text{CE}} &= -\mathbf{t}^{\top} (\log \mathbf{y}) \end{aligned}$$

Gradient descent updates can be derived for each row of W:

$$\frac{\partial L_{\text{CE}}}{\partial \mathbf{w}_k} = \frac{\partial L_{\text{CE}}}{\partial z_k} \cdot \frac{\partial z_k}{\partial \mathbf{w}_k} = (y_k - t_k) \cdot \mathbf{x}$$
$$\mathbf{w}_k \leftarrow \mathbf{w}_k - \alpha \frac{1}{N} \sum_{i=1}^{N} (y_k^{(i)} - t_k^{(i)}) \mathbf{x}^{(i)}$$

• Similar to linear/logistic reg (no coincidence) (verify the update)

#### What about SVMs?

- Not trivial to generalize the notion of a margin to multiclass setting.
- Many different proposals for multi-class SVMs, but outside of the scope of this course.

STA314-Lec6 29 / 44

- Let's imagine we have a prediction function  $y(\mathbf{x}, \theta)$  with parameters  $\theta$ , e.g.  $\theta = \mathbf{w}$  in logisitic regression.
- So far, the cost function  $\hat{\mathcal{R}}$  has been the average loss over the training examples:

$$\hat{\mathcal{R}}(\theta) = \frac{1}{N} \sum_{i=1}^{N} L^{(i)} = \frac{1}{N} \sum_{i=1}^{N} L(y(\mathbf{x}^{(i)}, \theta), t^{(i)}).$$

By linearity,

$$\frac{\partial \hat{\mathcal{R}}}{\partial \boldsymbol{\theta}} = \frac{1}{N} \sum_{i=1}^{N} \frac{\partial L^{(i)}}{\partial \boldsymbol{\theta}}.$$

• Computing the gradient requires summing over *all* of the training examples. This is known as batch training.

#### Stochastic Gradient Descent

- Batch training is impractical if you have a large dataset (e.g. millions of training examples)!
- Stochastic gradient descent (SGD): update the parameters based on the gradient for a single training example,

$$2- \qquad \boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \frac{\partial L^{(i)}}{\partial \boldsymbol{\theta}}$$

- Cost of each SGD update is independent of N!
- SGD can make significant progress before even seeing all the data!
- Mathematical justification: if you sample a training example uniformly at random, the stochastic gradient is an unbiased estimate of the batch gradient:

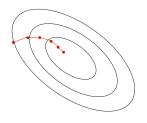
$$\mathbb{E}\left[\frac{\partial L^{(i)}}{\partial \boldsymbol{\theta}}\right] = \frac{1}{N} \sum_{i=1}^{N} \frac{\partial L^{(i)}}{\partial \boldsymbol{\theta}} = \frac{\partial \hat{\mathcal{R}}}{\partial \boldsymbol{\theta}}.$$

#### Stochastic Gradient Descent

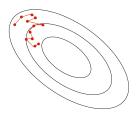
- Problems with using single training example to estimate gradient:
  - Variance in the estimate may be high
  - ▶ We can't exploit efficient vectorized operations
- Compromise approach:
  - ▶ compute the gradients on a randomly chosen medium-sized set of training examples  $\mathcal{M} \subset \{1, ..., N\}$ , called a mini-batch.
- Stochastic gradients computed on larger mini-batches have smaller variance.
- ullet The mini-batch size  $|\mathcal{M}|$  is a hyperparameter that needs to be set.
  - ▶ Too large: requires more compute; e.g., it takes more memory to store the activations, and longer to compute each gradient update
  - ▶ Too small: can't exploit vectorization, has high variance
  - ▶ A reasonable value might be  $|\mathcal{M}| = 100$ .

#### Stochastic Gradient Descent

• Batch gradient descent moves directly downhill. SGD takes steps in a noisy direction, but moves downhill on average.



batch gradient descent



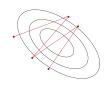
stochastic gradient descent

## Learning Rate

• In gradient descent, the learning rate  $\alpha$  is a hyperparameter we need to tune. Here are some things that can go wrong:



 $\alpha$  too small: slow progress



 $\alpha$  too large: oscillations

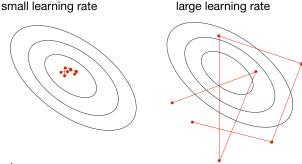


 $\alpha$  much too large: instability

• Good values are typically between 0.001 and 0.1. You should do a grid search if you want good performance (i.e. try 0.1, 0.03, 0.01, ...).

## SGD Learning Rate

 In stochastic training, the learning rate also influences the fluctuations due to the stochasticity of the gradients.



- Typical strategy:
  - Use a large learning rate early in training so you can get close to the optimum
  - ▶ Gradually decay the learning rate to reduce the fluctuations

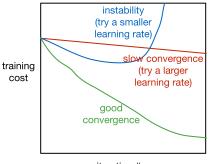
## SGD Learning Rate

 Warning: by reducing the learning rate, you reduce the fluctuations, which can appear to make the loss drop suddenly. But this can come at the expense of long-run performance.



## Training Curves

 To diagnose optimization problems, it's useful to look at training curves: plot the training cost as a function of iteration.



iteration #

 Warning: it's very hard to tell from the training curves whether an optimizer has converged. They can reveal major problems, but they can't guarantee convergence.

## The End of Supervised Learning

 This is the end of supervised learning in this course (sort of), next week we will move on to unsupervised learning. The following slides are optional and will not be tested.

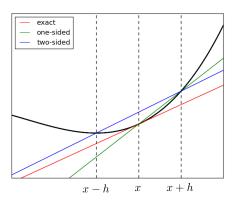
- We've derived a lot of gradients so far. How do we know if they're correct?
- Recall the definition of the partial derivative:

$$\frac{\partial}{\partial x_i} f(x_1, \dots, x_N) = \lim_{h \to 0} \frac{f(x_1, \dots, x_i + h, \dots, x_N) - f(x_1, \dots, x_i, \dots, x_N)}{h}$$

• Check your derivatives numerically by plugging in a small value of h, e.g.  $10^{-10}$ . This is known as finite differences.

• Even better: the two-sided definition

$$\frac{\partial}{\partial x_i} f(x_1, \dots, x_N) = \lim_{h \to 0} \frac{f(x_1, \dots, x_i + h, \dots, x_N) - f(x_1, \dots, x_i - h, \dots, x_N)}{2h}$$



- Run gradient checks on small, randomly chosen inputs
- Use double precision floats (not the default for TensorFlow, PyTorch, etc.!)
- Compute the relative error:

$$\frac{|a-b|}{|a|+|b|}$$

• The relative error should be very small, e.g.  $10^{-6}$ 

- Gradient checking is really important!
- Learning algorithms often appear to work even if the math is wrong.
- But:
  - ▶ They might work much better if the derivatives are correct.
  - Wrong derivatives might lead you on a wild goose chase.
- If you implement derivatives by hand, gradient checking is the single most important thing you need to do to get your algorithm to work well.