## CSC2515 Winter 2015 Introduction to Machine Learning

## Lecture 2: Linear regression

All lecture slides will be available as .pdf on the course website: <u>http://www.cs.toronto.edu/~urtasun/courses/CSC2515/CSC2515\_Winter15.html</u>

Many of the figures are provided by Chris Bishop from his textbook: "Pattern Recognition and Machine Learning"

## **Admin Details**

- Permanent tutorial time/place:
  - Tuesday 11-12, GB248
- Do I have the appropriate background?
  - Linear algebra: vector/matrix manipulations, properties
  - Calculus: partial derivatives
  - Probability: common distributions; Bayes Rule
  - Statistics: mean/median/mode; maximum likelihood
- Am I interested in this course?
  Many people in the waiting list
- Talk about Projects

## Outline

- Linear regression problem
  - continuous outputs
  - simple model
- Introduce key concepts:
  - loss functions
  - generalization
  - optimization
  - model complexity
  - regularization

### A simple example: 1-D regression

The green curve is the true function (which is not a polynomial) - not known

The data points are uniform in x but have noise in t.

$$f(x) = f(x) + \varepsilon$$

Aim: fit a curve to these points

from Bishop



Key questions:

- How do we parametrize the model (the curve)?
- What loss (objective) function should we use to judge fit?
- How do we optimize fit to unseen test data (generaliz., prediction)?

## **Example: Boston Housing data**

- Estimate median house price in a neighborhood based on neighborhood statistics
- Look at first (of 13) attributes: per capita crime rate
- Use this to predict house prices in other neighborhoods



#### Represent the Data

Data described as pairs D =  $((x^{(1)}, t^{(1)}), (x^{(2)}, t^{(2)}), ..., (x^{(N)}, t^{(N)}))$ 

- x is the input feature (per capita crime rate)
- t is the target output (median house price)
- Here t is continuous, so this is a regression problem

Could take first 300 examples as training set, remaining 206 as test set

- Use the training examples to construct hypothesis, or function approximator, that maps x to predicted y
- Evaluate hypothesis on test set

## Noise

A simple model typically does not exactly fit the data – lack of fit can be considered noise

Sources of noise

- Imprecision in data attributes (input noise)
- Errors in data targets (mislabeling)
- Additional attributes not taken into account by data attributes, affect target values (latent variables)
- Model may be too simple to account for data targets

#### Least-squares Regression

 Standard loss/cost/objective function measures the squared error in the prediction of t(x) from x.

$$y(x) = w_0 + w x$$

$$J(\mathbf{w}) = \sum_{n=1}^{N} [t^{(n)} - y(x^{(n)})]^2$$

 The loss for the red hypothesis is the sum of the squared vertical errors.



#### **Optimizing the Objective**

 One straightforward method: initialize w randomly, repeatedly update based on gradient descent in J

$$w \leftarrow w - \lambda \frac{\partial J}{\partial w}$$

- Here  $\lambda$  is the learning rate
- For a single training case, this gives the LMS update rule:

$$w \leftarrow w + 2\lambda(t^{(n)} - y(x^{(n)}))x^{(n)}$$

- Note: as error approaches zero, so does update
- What about w\_0?

## **Optimizing Across Training Set**

Two ways to generalize this for all examples in training set:

- 1. Stochastic/online updates update the parameters for each training case in turn, according to its own gradients
- 2. Batch updates: sum or average updates across every example i, then change the parameter values

$$w \leftarrow w + 2\lambda \sum_{n=1}^{N} (t^{(n)} - y(x^{(n)})) x^{(n)}$$

Underlying assumption: sample is independent and identically distributed (i.i.d.)

#### Non-iterative Least-squares Regression

An alternative optimization approach is non-iterative: take derivatives, set to zero, and solve for parameters.

$$\frac{dJ(\mathbf{w})}{dw_0} = -2\sum_{n=1}^{N} [t^{(n)} - (w_0 + w x^{(n)})] = 0$$
  
$$w_0 = (\sum_{n=1}^{N} t^{(n)} - w x^{(n)}) / N = \overline{t} - w\overline{x}$$
  
$$w = \frac{\sum_n (t^{(n)} - \overline{t})(x^{(n)} - \overline{x})}{\sum_n (x^{(n)} - \overline{x})^2}$$

11

#### When is minimizing the squared error equivalent to Maximum Likelihood Learning?

Minimizing the squared residuals is equivalent to maximizing the log probability of the correct answer under a Gaussian centered at the model's guess.

$$\mathcal{Y}^{(n)} = \mathcal{Y}(\mathbf{X}^{(n)}, \mathbf{W})$$

correct answer

t = the y = model'sestimate of most probable value

( )

$$p(t^{(n)} | y^{(n)}) = p(y^{(n)} + noise = t^{(n)} | \mathbf{x}^{(n)}, \mathbf{w}) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(t^{(n)} - y^{(n)})^2}{2\sigma^2}}$$
$$-\log p(t^{(n)} | y^{(n)}) = \log \sqrt{2\pi} + \log \sigma + \frac{(t^{(n)} - y^{(n)})^2}{2\sigma^2}$$
$$\text{Loss function}$$
$$\stackrel{\uparrow}{\underset{\text{can be ignored if sigma is fixed}}{\overset{\bullet}{\underset{\text{for every case}}}} can be ignored if$$

#### **Multi-dimensional inputs**

One method of extending the model is to consider other input dimensions

$$y(x) = w_0 + w_1 x_1 + w_2 x_2$$

In the Boston housing example, we can look at the number of rooms input feature

We can use gradient descent to solve for each coefficient, or use linear algebra - solve system of equations



## Linear Regression

- Imagine now want to predict the median house price from these multi-dimensional observations
- Each house is a data point n, with observations indexed by j:

$$\mathbf{x}^{(n)} = (x_1^{(n)}, ..., x_d^{(n)})$$

 Simple predictor is analogue of linear classifier, producing real-valued y for input x with parameters w (effectively fixing x<sub>0</sub> = 1):

$$\mathcal{Y} = \mathcal{W}_0 + \sum_{j=1}^d \mathcal{W}_j \mathcal{X}_j = \mathbf{W}^T \mathbf{X}$$

Why is the last equation correct?

14

## Linear models

- It is mathematically easy to fit linear models to data.
  - We can learn a lot about model-fitting in this relatively simple case.
- There are many ways to make linear models more powerful while retaining their nice mathematical properties:
  - By using non-linear, <u>non-adaptive</u> basis functions, we can get generalized linear models that learn non-linear mappings from input to output but are linear in their parameters - only the linear part of the model learns.
  - By using kernel methods we can handle expansions of the raw data that use a huge number of non-linear, non-adaptive basis functions.
- But linear methods will not solve most AI problems.
  - They have fundamental limitations.

#### Some types of basis functions in 1-D



Sigmoid and Gaussian basis functions can also be used in multilayer neural networks, but neural networks learn the parameters of the basis functions. This is much more powerful but also much harder and much messier. 16

## Two types of linear model that are equivalent with respect to learning $y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + w_2 x_2 + ... = \mathbf{w}^T \mathbf{x}$ $y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 \phi_1(\mathbf{x}) + w_2 \phi_2(\mathbf{x}) + ... = \mathbf{w}^T \Phi(\mathbf{x})$

- The first model has the same number of adaptive coefficients as the dimensionality of the data +1.
- The second model has the same number of adaptive coefficients as the number of basis functions +1.
- Once we have replaced the data by the outputs of the basis functions, fitting the second model is exactly the same problem as fitting the first model (unless we use the kernel trick)
- So we'll just focus on the first model

## Fitting a polynomial

 Now we use one of these basis functions: an M<sup>th</sup> order polynomial function

$$\mathcal{V}(x, \mathbf{w}) = \mathcal{W}_0 + \sum_{j=1}^M \mathcal{W}_j x^j$$

 We can use the same approaches to optimize the values of the weights on each coefficient: analytic, or iterative

## Minimizing squared error

Online Least mean squares: An alternative approach for really big datasets

$$\mathbf{w}^{(\tau+1)}_{\uparrow} = \mathbf{w}^{(\tau)} - \mathbf{\lambda} \nabla \mathbf{f}^{(n,\tau)}_{\uparrow}$$
weights after  
seeing training  
case tau+1 learning  
rate rate vector of derivatives of the  
squared error w.r.t. the  
weights on the training case  
presented at time tau.

- This is called "online" learning. It can be more efficient if the dataset is very redundant and it is simple to implement in hardware.
  - It is also called stochastic gradient descent if the training cases are picked at random.
  - Care must be taken with the learning rate to prevent divergent oscillations, and the rate must decrease at the end to get a good fit.

#### Some fits to the data: which is best?



#### 1-D regression illustrates key concepts

Data fits - is linear model best (model selection)?

- Simplest models do not capture all the important variations (signal) in the data: underfit
- More complex model may overfit the training data (fit not only the signal but also the noise in the data), especially if not enough data to constrain model

One method of assessing fit: test generalization = model's ability to predict the held out data

Optimization is essential: stochastic and batch iterative approaches; analytic when available

#### **Regularized least squares**

$$\tilde{\mathcal{J}}(\mathbf{w}) = \sum_{n=1}^{N} \{ \mathcal{V}(\mathbf{x}^{(n)}, \mathbf{w}) - t^{(n)} \}^2 + \alpha \mathbf{w}^T \mathbf{w}$$

The penalty on the squared weights (aka ridge regression, weight decay) is mathematically compatible with the squared error function, so we get a nice closed form for the optimal weights with this regularizer:

$$\mathbf{w}^* = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$$

$$\mathbf{\uparrow}$$
identity
matrix

Probabilistic Interpretation of the Regularizer

## A picture of the effect of the regularizer



- The overall cost function is the sum of two parabolic bowls.
- The sum is also a parabolic bowl.
- The combined minimum lies on the line between the minimum of the squared error and the origin.
- $w_1$  The regularizer just shrinks the weights.

#### A problem with the regularizer

- We would like the solution to be independent of the units we use to measure the components of the input vector.
- If different components have different units (e.g. age and height), we have a problem.
  - If we measure age in months and height in meters, the relative values of the two weights are very different than if we use years and millimeters  $\rightarrow$  squared penalty has very different effects
- One way to avoid the units problem: Whiten the data so that the input components all have unit variance and no covariance. This stops the regularizer from being applied to the whitening matrix.

$$\mathbf{X}_{whitened}^{T} = (\mathbf{X}^{T} \mathbf{X})^{-1/2} \mathbf{X}^{T}$$

- But this can cause other problems when two input components are almost perfectly correlated.
- We really need a prior on the weight on each input component, 25

## Other regularizers

- We do not need to use the squared error, provided we are willing to do more computation.
- Other powers of the weights can be used.



$$\tilde{J}(\mathbf{w}) = \sum_{n=1}^{N} \{ y(\mathbf{x}^{(n)}, \mathbf{w}) - t^{(n)} \}^2 + \alpha \| \mathbf{w} \|_q^q$$
<sup>26</sup>

# The lasso: penalizing the absolute values of the weights

$$\widetilde{\mathcal{J}}(\mathbf{w}) = \sum_{n=1}^{N} \left\{ \mathcal{V}(\mathbf{x}^{(n)}, \mathbf{w}) - t^{(n)} \right\}^2 + \alpha \sum_{j=1}^{d} |\mathcal{W}_j|$$

- Finding the minimum requires quadratic programming but its still unique because the cost function is convex (a bowl plus an inverted pyramid)
- As alpha is increased, many of the weights go to exactly zero.
  - This is great for interpretation, and it is also pretty good for preventing overfitting.

#### A geometrical view of the lasso compared with a penalty on the squared weights





Notice that  $w_1 = 0$  at the optimum

### Minimizing the absolute error

$$\min_{over \mathbf{w}} \sum_{n} |t^{(n)} - \mathbf{w}^T \mathbf{x}^{(n)}|$$

- This minimization involves solving a linear programming problem.
- It corresponds to maximum likelihood estimation if the output noise is modeled by a Laplacian instead of a Gaussian.

$$p(t^{(n)} | y^{(n)}) = \alpha e^{-\alpha |t^{(n)} - y^{(n)}|}$$
  
- log  $p(t^{(n)} | y^{(n)}) = -\alpha | t^{(n)} - y^{(n)} | + const$ 

#### One dimensional cross-sections of loss functions with different powers



#### The bias-variance decomposition



training datasets, of the model's estimate.

a derivation using a different notation (Bishop, 149)

# How the regularization parameter affects the bias and variance terms

high variance

low variance



#### An example of the bias-variance trade-off

