Machine Learning

Topic: Linear Regression Models

(contains ideas and a few images from wikipedia and books by Alpaydin, Duda/Hart/Stork, and Bishop. Updated Fall 2015)

Regression Learning Task

There is a set of possible examples $X = \{\mathbf{X}_1, \dots, \mathbf{X}_n\}$

Each example is a **vector** of k **real valued attributes**

$$\mathbf{x}_i = < x_{i1}, ..., x_{ik} >$$

There is a target function that maps X onto some real value Y $f: X \rightarrow Y$

The DATA is a set of tuples <example, response value>

$$\{\langle \mathbf{x}_1, y_1 \rangle, \dots \langle \mathbf{x}_n, y_n \rangle\}$$

Find a hypothesis **h** such that...

$$\forall \mathbf{x}, h(\mathbf{x}) \approx f(\mathbf{x})$$

Why use a linear regression model?

- Easily understood
- Interpretable
- Well studied by statisticians
 many variations and diagnostic measures
- Computationally efficient

Linear Regression Model

Assumption: The observed response (dependent) variable, y, is the true function, f(x), with additive Gaussian noise, ε , which has a mean of 0.

Deserved response
$$y = f(\mathbf{x}) + \mathcal{E}$$

Where $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

Assumption: The expected value of the response variable **y** is a linear combination of the k independent attributes/features)

The Hypothesis Space

Given the assumptions on the previous slide, our hypothesis space is the set of linear functions (hyperplanes)

$$h(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_k x_k$$

 $(w_0 \text{ is the offset from the origin. You always need } w_0)$

The goal is to learn a k+1 dimensional vector of weights that define a hyperplane minimizing an error criterion.

$$w = < w_0, w_1, ..., w_k >$$

Simple Linear Regression

- x has 1 attribute a (predictor variable)
- Hypothesis function is a line:



The Error Criterion

Typically estimate parameters by minimizing sum of squared residuals (RSS)...also known as the Sum of Squared Errors (SSE)



Multiple (Multivariate*) Linear Regression

- Many attributes $X_1, \ldots X_k$
- h(**x**) function is a hyperplane

*NOTE: In statistical literature, multivariate linear regression is regression with multiple outputs, and the case of multiple input variables is simply "multiple linear regression"



Formatting the data

Create a new 0 dimension with 1 and append it to the beginning of every example vector \mathbf{X}_i This placeholder corresponds to the offset \mathcal{W}_0

$$\mathbf{x}_i = <1, x_{i,1}, x_{i,2}, ..., x_{i,k} >$$

Format the data as a matrix of examples \mathbf{x} and a vector of response values y...



There is a closed-form solution!

Our goal is to find the weights of a function....

$$h(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_k x_k$$

...that minimizes the sum of squared residuals:

$$RSS = \sum_{i}^{n} (y_i - h(\mathbf{x}_i))^2$$

ท

It turns out that there is a close-form solution to this problem!

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Just plug your training data into the above formula and the best hyperplane comes out!

RSS in vector/matrix notation

$$RSS(\mathbf{w}) = \sum_{i=1}^{n} (y_i - h(\mathbf{x}_i))^2$$
$$= \sum_{i=1}^{n} (y_i - w_0 - \sum_{j=1}^{k} x_{ij} w_j)^2$$
$$= (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

Deriving the formula to find w

$$RSS(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^{T} (\mathbf{y} - \mathbf{X}\mathbf{w})$$
$$\frac{\partial RSS}{\partial \mathbf{w}} = -2\mathbf{X}^{T} (\mathbf{y} - \mathbf{X}\mathbf{w})$$
$$0 = -2\mathbf{X}^{T} (\mathbf{y} - \mathbf{X}\mathbf{w})$$
$$0 = \mathbf{X}^{T} (\mathbf{y} - \mathbf{X}\mathbf{w})$$
$$0 = \mathbf{X}^{T} (\mathbf{y} - \mathbf{X}\mathbf{w})$$
$$0 = \mathbf{X}^{T} \mathbf{y} - \mathbf{X}^{T} \mathbf{X}\mathbf{w}$$
$$\mathbf{X}^{T} \mathbf{X}\mathbf{w} = \mathbf{X}^{T} \mathbf{y}$$
$$\mathbf{w} = (\mathbf{X}^{T} \mathbf{X})^{-1} \mathbf{X}^{T} \mathbf{y}$$

What if X is not invertible?

• We said there was a closed form solution:

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- This presupposes matrix $(\mathbf{X}^T \mathbf{X})$ is invertible (non singular) and we can therefore find $(\mathbf{X}^T \mathbf{X})^{-1}$
- If two columns of X are exactly linearly related and thus not independent, then $({\bf X}^{ \mathrm{\scriptscriptstyle T}} {\bf X})$ is NOT invertible
- What then?

Your Friend: Dimensionality Reduction

- We need to make every column of X independent.
- The easy way: add a small amount random noise (with an expected value of 0) to X.
 - This is useful when you can't get rid of redundant columns for some reason
 - For example, your input data file is a 1000 examples of a constant value. You still want the code to return something, so you add a touch of noise and it will run and return something.
- The (often) better way: do dimensionality reduction to get rid of those redundant columns.

Making polynomial regression

You're familiar with linear regression where the input has k dimensions.

$$h(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_k x_k$$

We can use this same machinery to make polynomial regression from a one-dimensional input.....

$$h(x) = w_0 + w_1 x + w_2 x^2 + \dots + w_k x^k$$

Making polynomial regression

Given a scalar example z. We can make a k+1 dimensional example x $\mathbf{x} = \langle z^0, z^1, z^2, ..., z^k \rangle$

The *i*th element of x is the power z^i

$$h(x) = w_0 + w_1 z + w_2 z^2 + \dots w_k z^k$$

Making polynomial regression

Since $X_k \equiv z^k$ we can interpret the output of the regression as a polynomial function of Z

$$h(x) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_k x_k$$
$$= w_0 + w_1 z + w_2 z^2 + \dots + w_k z^k$$

Polynomial Regression

 Model the relationship between the response variable and the attributes/predictor variables as a kth-order polynomial. While this can model non-linear functions, it is still linear with respect to the coefficients.



Polynomial Regression

Parameter estimation (analytically minimizing sum of squared residuals):



(Note, there is only 1 attribute z for each training example. Those superscripts are powers, since we're doing polynomial regression)

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$







What happens if we fit to more data?

What happens if we fit to more data?

Bias and Variance of an Estimator

- Let X be a sample from a population specified by a true parameter θ
- Let d=d(X) be an estimator for θ

$$\mathbb{E}[(d-\theta)^2] = \mathbb{E}[(d-\mathbb{E}[d])^2] + (\mathbb{E}[d]-\theta)^2$$

mean square error

variance

bias²

As we **increase complexity**, **bias decreases** (a better fit to data) and **variance increases** (fit varies more with data)

Bias and Variance of Hypothesis Fn

• Bias:

Measures how much h(x) is wrong disregarding the effect of varying samples (This the statistical bias of an estimator. This is NOT the same as inductive bias, which is the set of assumptions that your learner is making)

• Variance:

Measures how much h(x) fluctuate around the expected value as the sample varies.

NOTE: These concepts are general machine learning concepts, not specific to linear regression.

Coefficient of Determination

the coefficient of determination, or R² indicates how well data points fit a line or curve. We'd like R² to be close to 1

$$R^{2} = 1 - E_{RSS} s$$
$$E_{RSS} = \frac{\sum_{i}^{n} (y_{i} - h(\mathbf{x}_{i}))^{2}}{\sum_{i}^{n} (y_{i} - \overline{y})^{2}} \text{ where } \overline{y} \text{ is the sample mean}$$

 $\mathbf{\cap}$

Don't just rely on numbers, visualize!

For all 4 sets: same mean and variance for x, same mean and variance (almost) for y, and same regression line and correlation between x and y (and therefore same R-squared).

Summary of Linear Regression Models

- Easily understood
- Interpretable
- Well studied by statisticians
- Computationally efficient
- Can handle non-linear situations if formulated properly
- Bias/variance tradeoff (occurs in all machine learning)
- Visualize!!
- GLMs

Appendix

(Stuff I couldn't cover in class)

high bias, low variance

high bias, high variance

low bias, high variance

low bias, low variance

• Bias:

Measures how much h(x) is wrong disregarding the effect of varying samples

• Variance:

Measures how much h(x) fluctuate around the expected value as the sample varies.

high variance → overfitting

There's a trade-off between bias and variance

Ways to Avoid Overfitting

- Simpler model
 - E.g. fewer parameters
- Regularization
 - penalize for complexity in objective function
- Fewer features
- Dimensionality reduction of features (e.g. PCA)
- More data...

Model Selection

- **Cross-validation**: Measure generalization accuracy by testing on data unused during training
- Regularization: Penalize complex models E'=error on data + λ model complexity Akaike's information criterion (AIC), Bayesian information criterion (BIC)
- Minimum description length (MDL): Kolmogorov complexity, shortest description of data
- Structural risk minimization (SRM)

Generalized Linear Models

- Models shown have assumed that the response variable follows a Gaussian distribution around the mean
- Can be generalized to response variables that take on *any* exponential family distribution (Generalized Linear Models - GLMs)