CPSC 340: Machine Learning and Data Mining

Nonlinear Regression Fall 2019

Last Time: Linear Regression

• We discussed linear models:

$$Y_{i} = w_{i} x_{i1} + w_{2} x_{i2} + \dots + w_{d} x_{id}$$

= $\sum_{s=1}^{d} w_{s} x_{ij} = w^{T} x_{i}$

- "Multiply feature x_{ij} by weight w_j, add them to get y_i".
- We discussed squared error function: $f(w) = \frac{1}{a} \sum_{i=1}^{n} (w^{T}x_{i} - y_{i})^{2}$ Predicted value
- Interactive demo:
 - http://setosa.io/ev/ordinary-least-squares-regression



To predict on test case \tilde{X}_i use $\tilde{y}_i = w^T \tilde{x}_i$

Matrix/Norm Notation (MEMORIZE/STUDY THIS)

- To solve the d-dimensional least squares, we use matrix notation:
 - We use 'w' as a "d times 1" vector containing weight ' w_i ' in position 'j'.
 - We use 'y' as an "n times 1" vector containing target ' y_i ' in position 'i'.
 - We use ' x_i ' as a "d times 1" vector containing features 'j' of example 'i'.
 - We're now going to be careful to make sure these are column vectors.
 - So 'X' is a matrix with x_i^T in row 'i'.



Matrix/Norm Notation (MEMORIZE/STUDY THIS)

- To solve the d-dimensional least squares, we use matrix notation:
 - Our prediction for example 'i' is given by the scalar $w^T x_i$.
 - Our predictions for all 'i' (n times 1 vector) is the matrix-vector product Xw.

$$y_{i} = w^{T}x_{i}$$

$$X_{w} = \begin{pmatrix} x_{i}^{T} \\ x_{2}^{T} \\ x_{1}^{T} \end{pmatrix} \begin{bmatrix} y \\ y_{2} \\ y_{2} \\ y_{3} \\ y_{4} \end{bmatrix} = \begin{pmatrix} x_{i}^{T}w \\ y_{2} \\ y_{4} \\ y_{4} \\ y_{4} \end{bmatrix} = \begin{pmatrix} x_{i}^{T}w \\ y_{2} \\ y_{4} \\ y_{4} \\ y_{4} \\ y_{4} \end{bmatrix} = \begin{pmatrix} x_{i}^{T}w \\ y_{4} \\$$

Matrix/Norm Notation (MEMORIZE/STUDY THIS)

- To solve the d-dimensional least squares, we use matrix notation:
 - Our prediction for example 'i' is given by the scalar $w^T x_i$.
 - Our predictions for all 'i' (n times 1 vector) is the matrix-vector product Xw.
 - Residual vector 'r' gives difference between predictions and y_i (n times 1).
 - Least squares can be written as the squared L2-norm of the residual.

Back to Deriving Least Squares for d > 2...

• We can write vector of predictions \hat{y}_i as a matrix-vector product:

$$\hat{\mathbf{y}} = \mathbf{x}_{\mathbf{w}} = \begin{pmatrix} \mathbf{w}_{\mathbf{x}_{1}} \\ \mathbf{w}_{\mathbf{x}_{1}} \\ \mathbf{w}_{\mathbf{x}_{n}} \\ \mathbf{w}_{\mathbf{x}_{n}} \end{pmatrix}$$

• And we can write linear least squares in matrix notation as:

$$f(w) = \frac{1}{2} || x_w - y ||^2 = \frac{1}{2} \sum_{i=1}^{2} (w_{x_i} - y_i)^2$$

We'll use this notation to derive d-dimensional least squares 'w'.
 By setting the gradient ∇ f(w) equal to the zero vector and solving for 'w'.

Digression: Matrix Algebra Review

- Quick review of linear algebra operations we'll use:
 - If 'a' and 'b' be vectors, and 'A' and 'B' be matrices then:

$$a^{T}b = b^{T}a$$

$$\|a\|^{2} = a^{T}a$$

$$(A + B)^{T} = A^{T} + B^{T}$$

$$(AB)^{T} = B^{T}A^{T}$$

$$(A+B)(A+B) = AA + BA + AB + BB$$

$$a^{T}AL = b^{T}A^{T}a$$

$$\bigvee_{vector} \qquad \bigvee_{vector}$$

Sanity check: ALWAYS CHECK THAT DIMENSIONS MATCH (if not, you did something wrong)

Linear and Quadratic Gradients

• From these rules we have (see post-lecture slide for steps):

$$f(u) = \frac{1}{2} \sum_{i=1}^{2} (u^{T} x_{i}^{-} y_{i})^{2} = \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - w^{T} X^{T} x_{w} - w^{T} X^{T} y_{v} + \frac{1}{2} y^{T} y_{v}$$

$$= \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - w^{T} X^{T} x_{w} - w^{T} X^{T} y_{v} + \frac{1}{2} y^{T} y_{v}$$

$$= \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - w^{T} X^{T} x_{w} - w^{T} X^{T} y_{v} + \frac{1}{2} y^{T} y_{v}$$

$$= \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{$$

J

• How do we compute gradient?

Let's first do it with
$$d=1$$
:
 $f(w) = \frac{1}{2}waw + wb + c$
 $= \frac{1}{2}aw^{2} + wb + c$
 $f'(w) = aw + b+0$
 $f'(w) = aw + b+$

Linear and Quadratic Gradients

• We've written as a d-dimensional quadratic:

$$f(u) = \frac{1}{2} \sum_{i=1}^{2} (w^{T} x_{i}^{-} y_{i})^{2} = \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - w^{T} X^{T} u - w^{T} X^{T} y + \frac{1}{2} y^{T} y|$$

$$= \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - w^{T} X^{T} u - w^{T} X^{T} y + \frac{1}{2} y^{T} y|$$

$$= \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - w^{T} X^{T} u - w^{T} X^{T} y + \frac{1}{2} y^{T} y|$$

$$= \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - w^{T} X^{T} u - w^{T} X^{T} y + \frac{1}{2} y^{T} y|$$

$$= \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - w^{T} X^{T} u - w^{T} X^{T} y + \frac{1}{2} y^{T} y|$$

$$= \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - w^{T} X^{T} u - w^{T} X^{T} y + \frac{1}{2} y^{T} y|$$

$$= \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - w^{T} X^{T} u - w^{T} X^{T} y + \frac{1}{2} ||X_{w} - y||^{2} = \frac{1}{2} ||X_{w} - y||^{2} =$$

- Gradient is given by: $\nabla f(w) = Aw b + D$
- Using definitions of 'A' and 'b': $= \chi^T \chi_w \chi^T \gamma$

Sanity check: all dimensions match
$$(d \times n) (n \times d) (d \times 1) - (d \times n) (n \times 1)$$

Normal Equations

- Set gradient equal to zero to find the "critical" points: $\chi^{\gamma}\chi_{u} - \chi^{\gamma}\gamma = O$
- We now move terms not involving 'w' to the other side:

$$\chi^{\gamma}\chi_{w} = \chi^{\gamma}\gamma$$

- This is a set of 'd' linear equations called the normal equations.
 - This a linear system like "Ax = b" from Math 152.
 - You can use Gaussian elimination to solve for 'w'.
 - In Julia, the "\" command can be used to solve linear systems:

Train:
$$W = (X'X) \setminus (X'y)$$
 Predict: yhat = $X_{lost} * W$

Incorrect Solutions to Least Squares Problem

The least synares objective is
$$F(w) = \frac{1}{2} ||Xw - y||^2$$

The minimizers of this objective are solutions to the linear system:
 $X^T X w = X^7 y$
The following are not the solutions to the least squares problem:
 $w = (X^T X)^{-1} (X^7 y)$ (only true if $X^T X$ is invertible)
 $w X^T X = X^7 y$ (matrix multiplication is not commutative, dimensions don'
 $w = \frac{X^T y}{X^T X}$ (you cannot divide by a matrix)

Least Squares Cost

- Cost of solving "normal equations" X^TXw = X^Ty?
- Forming X^Ty vector costs O(nd).

- It has 'd' elements, and each is an inner product between 'n' numbers.

• Forming matrix X^TX costs O(nd²).

- It has d² elements, and each is an inner product between 'n' numbers.

- Solving a d x d system of equations costs O(d³).
 - Cost of Gaussian elimination on a d-variable linear system.
 - Other standard methods have the same cost.
- Overall cost is O(nd² + d³).
 - Which term dominates depends on 'n' and 'd'.

Least Squares Issues

- Issues with least squares model:
 - Solution might not be unique.
 - It is sensitive to outliers.
 - It always uses all features.
 - Data can might so big we can't store $X^T X$.
 - Or you can't afford the O(nd² + d³) cost.
 - It might predict outside range of y_i values.
 - It assumes a linear relationship between x_i and y_i.

>X is nxd so XT is dxn and XTX is dxd.

Non-Uniqueness of Least Squares Solution

- Why isn't solution unique?
 - Imagine having two features that are identical for all examples.
 - I can increase weight on one feature, and decrease it on the other, without changing predictions.

$$\gamma_{i} = w_{1} \chi_{i1} + w_{2} \chi_{i1} = (w_{1} + w_{2}) \chi_{i1} + 0 \chi_{i1}$$

- Thus, if (w_1, w_2) is a solution then $(w_1+w_2, 0)$ is another solution.
- This is special case of features being "collinear":
 - One feature is a linear function of the others.
- But, any 'w' where $\nabla f(w) = 0$ is a global minimizer of 'f'.
 - This is due to convexity of 'f', which we'll discuss later.

(pause)

Motivation: Non-Linear Progressions in Athletics

• Are top athletes going faster, higher, and farther?



HIGH JUMP PROGRESSION MEN AND WOMEN (mean of top ten)











http://www.at-a-lanta.nl/weia/Progressie.html https://en.wikipedia.org/wiki/Usain_Bolt http://www.britannica.com/biography/Florence-Griffith-Joyner

• We can adapt our classification methods to perform regression:

- We can adapt our classification methods to perform regression:
 - Regression tree: tree with mean value or linear regression at leaves.



http://www.at-a-lanta.nl/weia/Progressie.html

- We can adapt our classification methods to perform regression:
 - Regression tree: tree with mean value or linear regression at leaves.
 - Probabilistic models: fit $p(x_i | y_i)$ and $p(y_i)$ with Gaussian or other model.
 - CPSC 540.



- We can adapt our classification methods to perform regression:
 - Regression tree: tree with mean value or linear regression at leaves.
 - Probabilistic models: fit $p(x_i | y_i)$ and $p(y_i)$ with Gaussian or other model.
 - Non-parametric models:
 - KNN regression:
 - Find 'k' nearest neighbours of \tilde{X}_{i} .
 - Return the mean of the corresponding y_i.



- We can adapt our classification methods to perform regression:
 - Regression tree: tree with mean value or linear regression at leaves.
 - Probabilistic models: fit $p(x_i | y_i)$ and $p(y_i)$ with Gaussian or other model.
 - Non-parametric models:
 - KNN regression.
 - Could be weighted by distance.
 - Close points 'j' get more "weight" w_{ij}.



- We can adapt our classification methods to perform regression:
 - Regression tree: tree with mean value or linear regression at leaves.
 - Probabilistic models: fit $p(x_i | y_i)$ and $p(y_i)$ with Gaussian or other model.
 - Non-parametric models:
 - KNN regression.
 - Could be weighted by distance.
 - 'Nadaraya-Waston': weight *all* y_i by distance to x_i.²⁵

$$\hat{y}_{i} = \frac{\sum_{j=1}^{n} v_{ij} y_{j}}{\sum_{j=1}^{n} v_{ij}}$$



http://www.mathworks.com/matlabcentral/fileexchange/35316-kernel-regression-with-variable-window-width/content/ksr_vw.m

Adapting Counting/

- We can adapt our classification
 - Regression tree: tree with mea >
 - Probabilistic models: fit $p(x_i | y$
 - Non-parametric models:
 - KNN regression.
 - Could be weighted by distance.
 - 'Nadaraya-Waston': weight *all* y_i
 - 'Locally linear regression': for each x_i, fit a linear model weighted by distance.

(Better than KNN and NW at boundaries.)



- We can adapt our classification methods to perform regression:
 - Regression tree: tree with mean value or linear regression at leaves.
 - Probabilistic models: fit $p(x_i | y_i)$ and $p(y_i)$ with Gaussian or other model.
 - Non-parametric models:
 - KNN regression.
 - Could be weighted by distance.
 - 'Nadaraya-Waston': weight *all* y_i by distance to x_i.
 - 'Locally linear regression': for each x_i, fit a linear model weighted by distance.
 (Better than KNN and NW at boundaries.)
 - Ensemble methods:
 - Can improve performance by averaging across regression models.

- We can adapt our classification methods to perform regression.
- Applications:
 - Regression forests for fluid simulation:
 - https://www.youtube.com/watch?v=kGB7Wd9CudA
 - KNN for image completion:
 - <u>http://graphics.cs.cmu.edu/projects/scene-completion</u>
 - Combined with "graph cuts" and "Poisson blending".
 - KNN regression for "voice photoshop":
 - https://www.youtube.com/watch?v=I3I4XLZ59iw
 - Combined with "dynamic time warping" and "Poisson blending".
- But we'll focus on linear models with non-linear transforms.
 - These are the building blocks for more advanced methods.

Why don't we have a y-intercept?

- Linear model is $\hat{y}_i = wx_i$ instead of $\hat{y}_i = wx_i + w_0$ with y-intercept w_0 .
- Without an intercept, if $x_i = 0$ then we must predict $\hat{y}_i = 0$.



Why don't we have a y-intercept?

- Linear model is $\hat{y}_i = wx_i$ instead of $\hat{y}_i = wx_i + w_0$ with y-intercept w_0 .
- Without an intercept, if $x_i = 0$ then we must predict $\hat{y}_i = 0$.



Adding

res this

Adding a Bias Variable

- Simple trick to add a y-intercept ("bias") variable:
 - Make a new matrix "Z" with an extra feature that is always "1".

$$X = \begin{bmatrix} -0.1 \\ 0.3 \\ 0.2 \end{bmatrix} \qquad \qquad Z = \begin{bmatrix} 1 & -0.1 \\ 1 & 0.3 \\ 1 & 0.2 \end{bmatrix}$$

- Now use "Z" as your features in linear regression.
 - We'll use 'v' instead of 'w' as regression weights when we use features 'Z'.

$$\dot{y}_{i} = V_{1} Z_{i1} + V_{2} Z_{i2} = W_{0} + W_{1} X_{i1}$$

$$\int_{W_{0}} \int_{W_{1}} \int_{W_{1}} \int_{W_{1}} \int_{X_{1}} \int_{W_{1}} X_{i1}$$

- So we can have a non-zero y-intercept by changing features.
 - This means we can ignore the y-intercept in our derivations, which is cleaner.

Motivation: Limitations of Linear Models

• On many datasets, y_i is not a linear function of x_i.



• Can we use least square to fit non-linear models?

Non-Linear Feature Transforms

- Can we use linear least squares to fit a quadratic model? $\hat{y_i} = w_{\partial} + w_i x_i + w_2 x_i^2$
- You can do this by changing the features (change of basis):

$$X = \begin{bmatrix} 6,2\\ -0.5\\ 1\\ 4 \end{bmatrix} \qquad Z = \begin{bmatrix} 1 & 0.2 & (0.2)^2\\ 1 & -0.5 & (-0.5)^2\\ 1 & 1 & (1)^2\\ 1 & 4 & (4)^2 \end{bmatrix}$$

$$Y^{-inf} X \qquad x^2$$

- Fit new parameters 'v' under "change of basis": solve $Z^TZv = Z^Ty$.
- It's a linear function of w, but a quadratic function of x_i.

$$\hat{y}_{i} = \hat{v}_{Z_{i}}^{T} = \hat{v}_{Z_{i}}^{T} + \hat{v}_{Z_{i2}}^{T} + \hat{v}_{Z_{i3}}^{T} + \hat{v}_{Z_{i3}}^{T}$$

Non-Linear Feature Transforms



General Polynomial Features (d=1)

• We can have a polynomial of degree 'p' by using these features:

$$Z = \begin{bmatrix} 1 & x_{1} & (x_{1})^{2} & \dots & (x_{n})^{p} \\ 1 & x_{2} & (x_{2})^{2} & \dots & (x_{n})^{p} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n} & (x_{n})^{2} & \dots & (x_{n})^{p} \end{bmatrix}$$

- There are polynomial basis functions that are numerically nicer:
 - E.g., Lagrange polynomials (see CPSC 303).

Summary

- Matrix notation for expressing least squares problem.
- Normal equations: solution of least squares as a linear system.
 Solve (X^TX)w = (X^Ty).
- Solution might not be unique because of collinearity.
 But any solution is optimal because of "convexity".
- Tree/probabilistic/non-parametric/ensemble regression methods.
- Non-linear transforms:
 - Allow us to model non-linear relationships with linear models.

• Next time: how to do least squares with a million features.

Linear Least Squares: Expansion Step

Wont 'u' that minimizes

$$f(w) = \frac{1}{2} \sum_{i=1}^{n} (w^{T}x_{i} - y_{i})^{2} = \frac{1}{2} ||Xw - y||_{2}^{2} = \frac{1}{2} (Xw - y)^{T} (Xw - y) \qquad ||a||^{2} = a^{T}a$$

$$= \frac{1}{2} ((xw)^{T} - y^{T}) (Xw - y) \qquad (A+b^{T}) = (A^{T}+b^{T})$$

$$= \frac{1}{2} (w^{T}X^{T} - y^{T}) (Xw - y) \qquad (Ab)^{T} = B^{T}A^{T}$$

$$= \frac{1}{2} (w^{T}X^{T} (Xw - y) - y^{T} (Xw - y)) (A+b)(=AC+bC)$$

$$= \frac{1}{2} (w^{T}X^{T} Xw - w^{T}X^{T}y - y^{T}Xw + y^{T}y) \qquad A(b+c)=AbbBC$$

$$= \frac{1}{2} (w^{T}X^{T} Xw - w^{T}X^{T}y + \frac{1}{2}y^{T}y) \qquad A(b+c)=Ab^{T}A^{T}a$$

$$= \frac{1}{2} w^{T}X^{T}Xw - w^{T}X^{T}y + \frac{1}{2}y^{T}y \qquad A(b+c)=Ab^{T}A^{T}a$$

$$= \frac{1}{2} w^{T}X^{T}Xw - w^{T}X^{T}y + \frac{1}{2}y^{T}y \qquad A(b+c)=Ab^{T}A^{T}a$$

$$= \frac{1}{2} w^{T}X^{T}Xw - w^{T}X^{T}y + \frac{1}{2}y^{T}y \qquad A(b+c)=Ab^{T}A^{T}a$$

$$= \frac{1}{2} w^{T}X^{T}Xw - w^{T}X^{T}y + \frac{1}{2}y^{T}y \qquad A(b+c)=Ab^{T}A^{T}a$$

Vector View of Least Squares

• We showed that least squares minimizes:

$$f(w) = \frac{1}{2} ||X_w - y||^2$$

- The ½ and the squaring don't change solution, so equivalent to: $f(w) = \|\chi_w - \gamma\|$
- From this viewpoint, least square minimizes Euclidean distance between vector of labels 'y' and vector of predictions Xw.

Bonus Slide: Householder(-ish) Notation

 Househoulder notation: set of (fairly-logical) conventions for math. Use greek letters for scalors &= 1, B= 35, 7= 11 Use <u>first/last lowercase</u> letters for vectors: $w = \begin{bmatrix} 0 & 1 \\ 0 & 2 \end{bmatrix}$, $x = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$, $y = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$, $a = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$, $b = \begin{bmatrix} 0 & 5 \\ 0 & 5 \end{bmatrix}$ Assumed to be column-vectors. Use first/last uppercase letters for matrices: X, Y, W, A, B Indices use i, j, K. Sizes use m, n, d, p, and k is obvious from context Sets use S, T, U, V When I write x, I Functions use f, g, and h. mean "grab row 'i' of X and make a column-vector with its values."

Bonus Slide: Householder(-ish) Notation

• Househoulder notation: set of (fairly-logical) conventions for math:

Our ultimate least squares notation:

$$f(w) = \frac{1}{2} ||Xw - y||^2$$

But if we agree on notation we can quickly understand:

$$g(x) = \frac{1}{2} ||Ax - b||^2$$

If we use random notation we get things like:

$$H(\beta) = \frac{1}{2} ||R\beta - P_n||^2$$
Is this the same mode

12

When does least squares have a unique solution?

- We said that least squares solution is not unique if we have repeated columns.
- But there are other ways it could be non-unique:
 - One column is a scaled version of another column.
 - One column could be the sum of 2 other columns.
 - One column could be three times one column minus four times another.
- Least squares solution is unique if and only if all columns of X are "linearly independent".
 - No column can be written as a "linear combination" of the others.
 - Many equivalent conditions (see Strang's linear algebra book):
 - X has "full column rank", $X^T X$ is invertible, $X^T X$ has non-zero eigenvalues, det($X^T X$) > 0.
 - Note that we cannot have independent columns if d > n.