# CPSC 340:
# Machine Learning and Data Mining

Nonlinear Regression

Fall 2020

# Last Time: Linear Regression

- We discussed linear models:

$$y_i = w_1 x_{i1} + w_2 x_{i2} + \cdots + w_d x_{id}$$
$$= \sum_{j=1}^{d} w_j x_{ij} = w^T x_i$$

- "Multiply feature $x_{ij}$ by weight $w_j$, add them to get $y_i$".
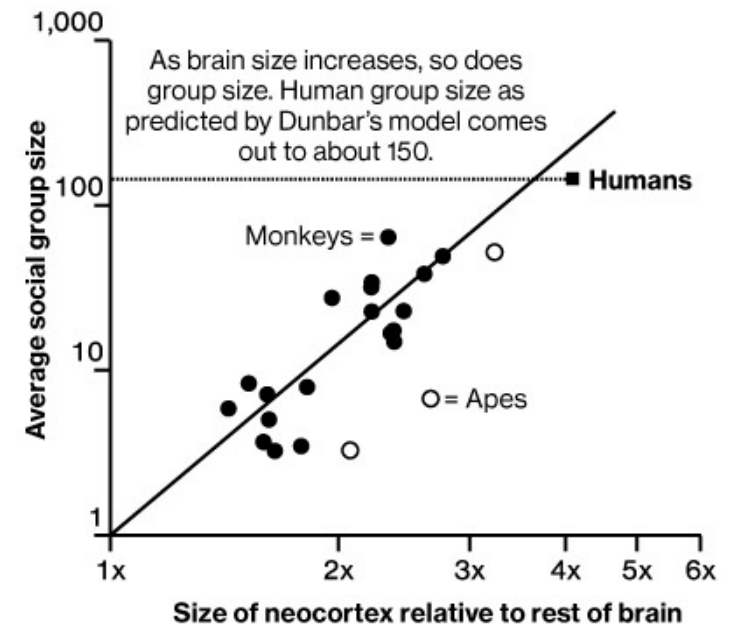
- We discussed squared error function:

$$f(w) = \frac{1}{2} \sum_{i=1}^{n} (w^T x_i - y_i)^2$$

Predicted value ⟵     ⟶ True value

- Interactive demo:
  - http://setosa.io/ev/ordinary-least-squares-regression

**The Social Cortex**

1,000

As brain size increases, so does group size. Human group size as predicted by Dunbar's model comes out to about 150.

100 ······· ■ **Humans**

Monkeys = ●

Average social group size

10

O = Apes

1

1x   2x   3x   4x   5x   6x

Size of neocortex relative to rest of brain

DATA: THE SOCIAL BRAIN HYPOTHESIS, DUNBAR 1998

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

# 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 '$w_j$' 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'.

$$
w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}
\quad
y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}
\quad
x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{bmatrix}
\quad
X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & & & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nd} \end{bmatrix}
=
\begin{bmatrix} - x_1^T - \\ - x_2^T - \\ \vdots \\ - x_n^T - \end{bmatrix}
$$

# 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 $Xw$.

$$\hat{y}_i = w^T x_i$$

Also, because $w^T x_i$ is a scalar,
we have $w^T x_i = x_i^T w$.

$$(e.g., [5]^T = [5])$$

$$Xw = \begin{bmatrix} - x_1^T - \\ - x_2^T - \\ \vdots \\ - x_n^T - \end{bmatrix} \begin{bmatrix} | \\ w \\ | \end{bmatrix} = \begin{bmatrix} x_1^T w \\ x_2^T w \\ \vdots \\ x_n^T w \end{bmatrix} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \hat{y}$$

$X$

$w$

Prediction for example 'i' 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 $Xw$.
  - 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.

$$f(w) = \sum_{i=1}^{n} (w^T x_i - y_i)^2 = \sum_{i=1}^{n} (r_i)^2$$

$$r_i$$

$$r = \hat{y} - y = Xw - y = \begin{bmatrix} w^T y_1 \\ w^T x_2 \\ \vdots \\ w^T x_n \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} w^T x_1 - y_1 \\ w^T x_2 - y_2 \\ \vdots \\ w^T x_n - y_n \end{bmatrix}$$

$$\hat{y}$$

$r_2$ is difference for example 2.

$$= \sum_{i=1}^{n} r_i r_i$$

$$= r^T r$$

$$= \|r\|^2 = \|Xw - y\|^2$$

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

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

$$\hat{y} = Xw = \begin{bmatrix} w^T x_1 \\ w^T x_2 \\ \vdots \\ w^T x_n \end{bmatrix}$$

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

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

- We'll use this notation to derive d-dimensional least squares 'w'.
  - By setting the gradient $\nabla f(w)$ 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:

$$a^T b = b^T a$$

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

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

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

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

$$a^T \underbrace{Ab}_{vector} = b^T \underbrace{A^T a}_{vector}$$

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):

$$f(w) = \frac{1}{2}\sum_{i=1}^{n}(w^\top x_i - y_i)^2 = \frac{1}{2}\|Xw - y\|^2 = \frac{1}{2}\underbrace{w^\top X^\top X w}_{\text{matrix 'A'}} - \underbrace{w^\top X^\top y}_{\text{vector 'b'}} + \underbrace{\frac{1}{2}y^\top y}_{\text{scalar 'c'}}$$

$$= \frac{1}{2}w^\top A w + w^\top b + c \quad \rightarrow \text{These are scalars so dimensions match.}$$

- How do we compute gradient?

Let's first do it with $d=1$:

$$f(w) = \frac{1}{2}waw + wb + c$$
$$= \frac{1}{2}aw^2 + wb + c$$

$$f'(w) = aw + b + 0$$

Here are the generalizations to 'd' dimensions:

$$\nabla[c] = 0 \quad \text{(zero vector)}$$
$$\nabla[w^\top b] = b$$
$$\nabla[\tfrac{1}{2}w^\top Aw] = Aw \quad \text{(if A is symmetric)}$$

→ Full derivations are on webpage in notes on linear and quadratic gradients.

# Linear and Quadratic Gradients

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

$$f(w) = \frac{1}{2} \sum_{i=1}^{n} (w^\top x_i - y_i)^2 = \frac{1}{2} \|Xw - y\|^2 = \frac{1}{2} \underbrace{w^\top X^\top X w}_{\text{matrix 'A'}} - \underbrace{w^\top X^\top y}_{\text{vector 'b'}} + \underbrace{\frac{1}{2} y^\top y}_{\text{scalar 'c'}}$$

$$= \frac{1}{2} w^\top A w + w^\top b + c$$

- Gradient is given by: $\quad \nabla f(w) = Aw - b + 0$

- Using definitions of 'A' and 'b': $\quad = X^\top X w - X^\top y$

  $\underbrace{\phantom{= X^\top X w}}$

  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:

$$X^TX_w - X^Ty = 0$$

- We now move terms not involving 'w' to the other side:

$$X^TX_w = X^Ty$$

- 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:

$$\text{Train: } w = (X'X) \backslash (X'y) \qquad \text{Predict: } yhat = X_{test} * w$$

# Normal Equations

- Set gradient equal to zero to find the "critical" points:

$$X^T X_w - X^T y = 0$$

- We now move terms not involving 'w' to the other side:

$$X^T X_w = X^T y$$

- 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 Python, you solve linear systems in 1 line using numpy.linalg.solve.

# Incorrect Solutions to Least Squares Problem

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

The minimizers of this objective are solutions to the linear system:

$$X^T X w = X^T y$$

The following are <u>not</u> the solutions to the least squares problem:

$w = (X^T X)^{-1} (X^T y)$  (only true if <u>$X^T X$ is invertible</u>)

$w X^T X = X^T y$  (matrix multiplication is <u>not</u> commutative, dimensions don't even match)

$w = \dfrac{X^T y}{X^T X}$  (you <u>cannot divide by a matrix</u>)

# 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 $O(nd^2)$.
  - It has $d^2$ elements, and each is an inner product between 'n' numbers.
- Solving a d x d system of equations costs $O(d^3)$.
  - Cost of Gaussian elimination on a d-variable linear system.
  - Other standard methods have the same cost.
- Overall cost is $O(nd^2 + d^3)$.
  - 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 $X^TX$.
    - Or you can't afford the $O(nd^2 + d^3)$ cost.
  - It might predict outside range of $y_i$ values.
  - It assumes a linear relationship between $x_i$ and $y_i$.

$X$ is $n \times d$

so $X^T$ is $d \times n$

and $X^TX$ is $d \times d$.

# 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.

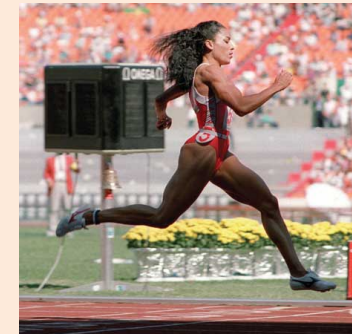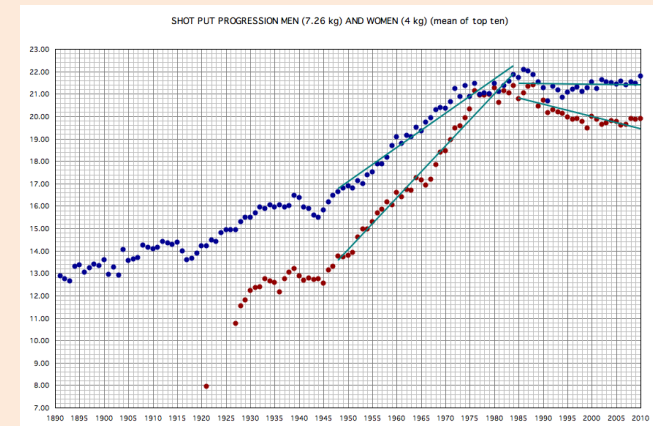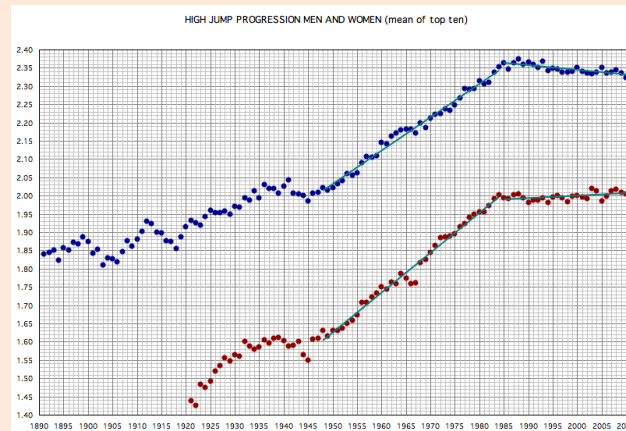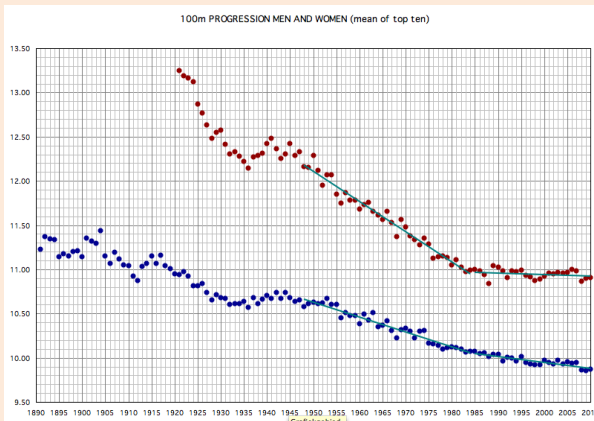$$\hat{y}_i = w_1 x_{i1} + w_2 x_{i1} = (w_1 + w_2)x_{i1} + 0 x_{i1}$$

  copy

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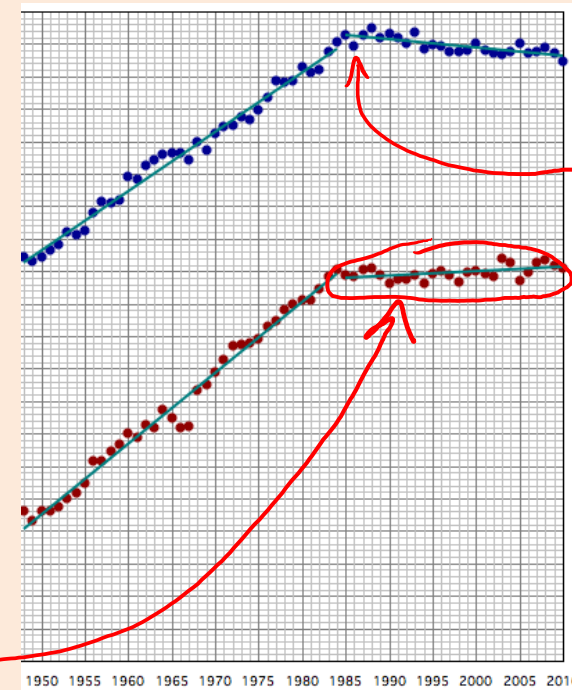
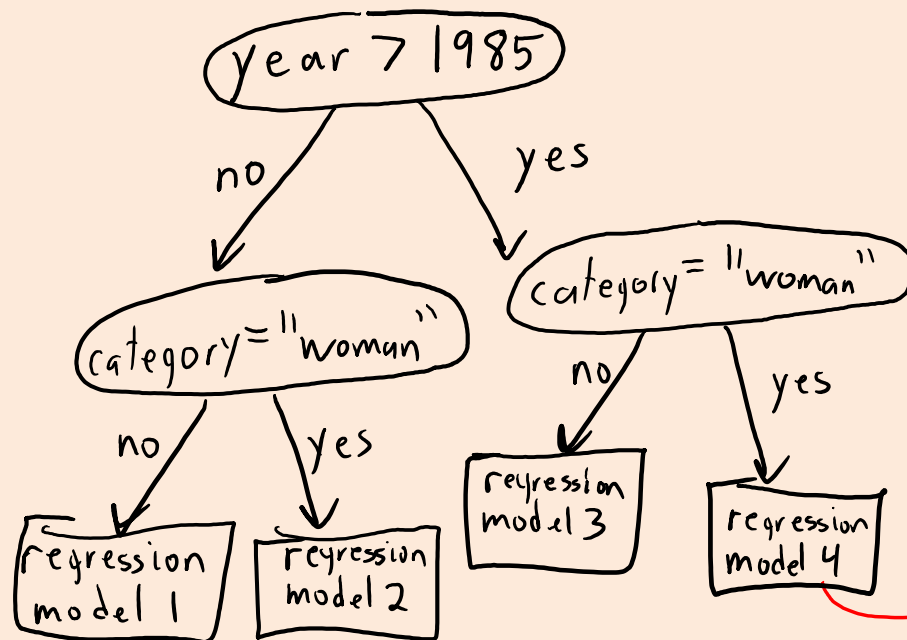- Are top athletes going faster, higher, and farther?

# Adapting Counting/Distance-Based Methods

- We can adapt our classification methods to perform regression:
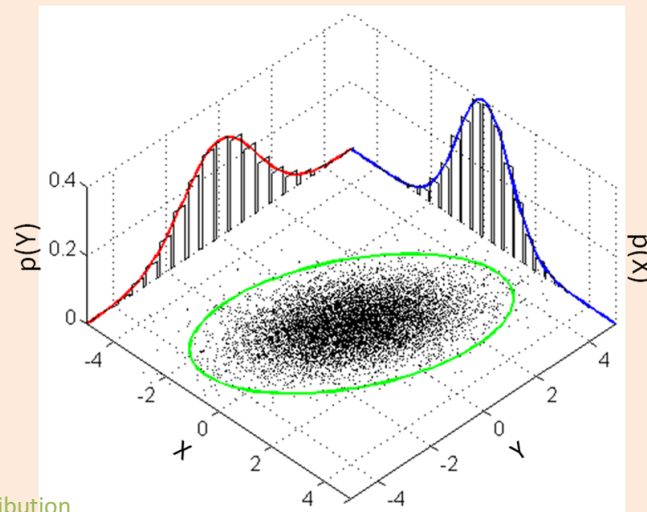
# Adapting Counting/Distance-Based Methods

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

# Adapting Counting/Distance-Based Methods

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

# Adapting Counting/Distance-Based Methods

- We can adapt our classification methods to perform regression:
  - Regression tree: tree with mean value or linear regression at leaves.
  - Probabilistic models: fit $p(x_i \mid y_i)$ 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 $y_i$.
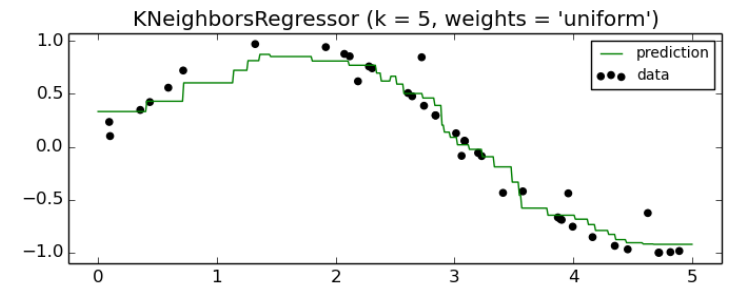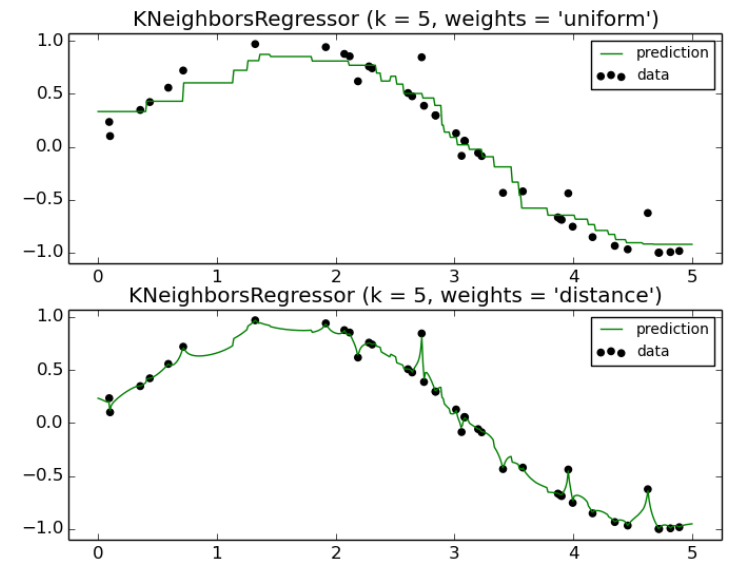


KNeighborsRegressor (k = 5, weights = 'uniform')

# Adapting Counting/Distance-Based Methods

- We can adapt our classification methods to perform regression:
    - Regression tree: tree with mean value or linear regression at leaves.
    - Probabilistic models: fit $p(x_i \mid y_i)$ 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" $w_{ij}$.
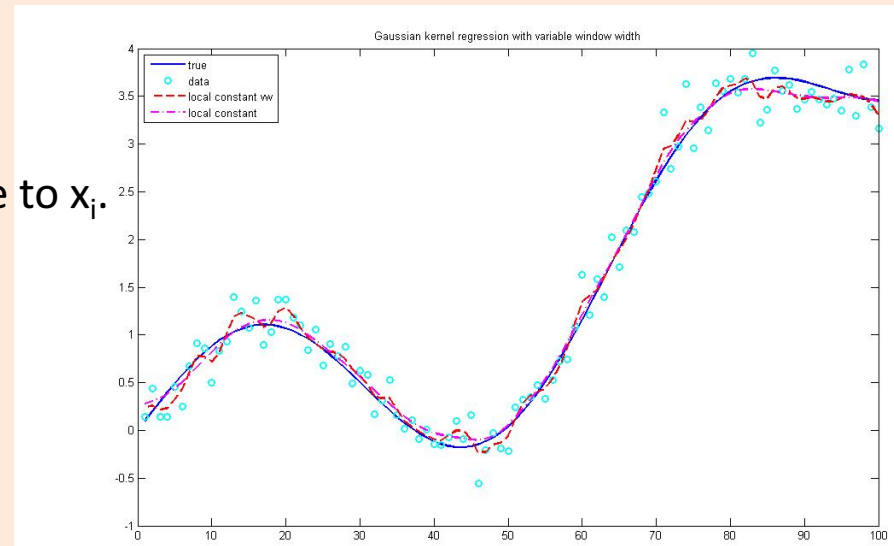
# Adapting Counting/Distance-Based Methods

- We can adapt our classification methods to perform regression:
  - Regression tree: tree with mean value or linear regression at leaves.
  - Probabilistic models: fit $p(x_i \mid y_i)$ 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 $x_i$.

$$\hat{y}_i = \frac{\sum_{j=1}^{n} v_{ij} y_j}{\sum_{j=1}^{n} v_{ij}}$$



Gaussian kernel regression with variable window width

true
data
local constant vw
local constant

# Adapting Counting/

- We can adapt our classific[...]
  - Regression tree: tree with mea[...]
  - Probabilistic models: fit $p(x_i \mid y$[...]
  - 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.)



d=2, q=0.5

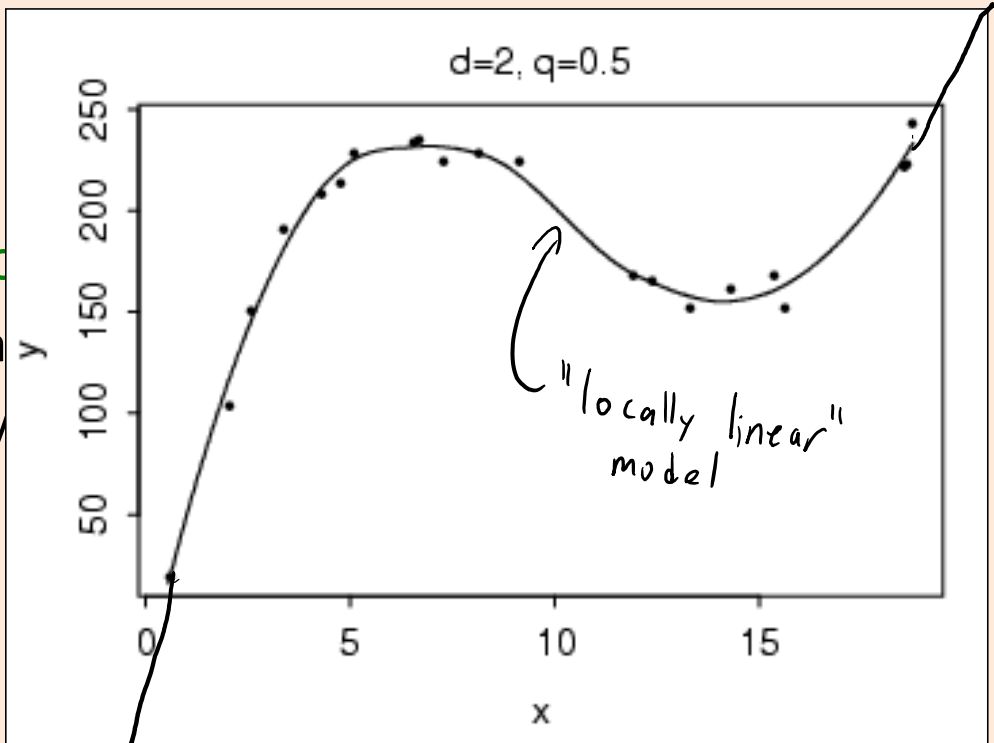"locally linear" model

# Adapting Counting/Distance-Based Methods

- We can adapt our classification methods to perform regression:
  - Regression tree: tree with mean value or linear regression at leaves.
  - Probabilistic models: fit $p(x_i \mid y_i)$ 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.
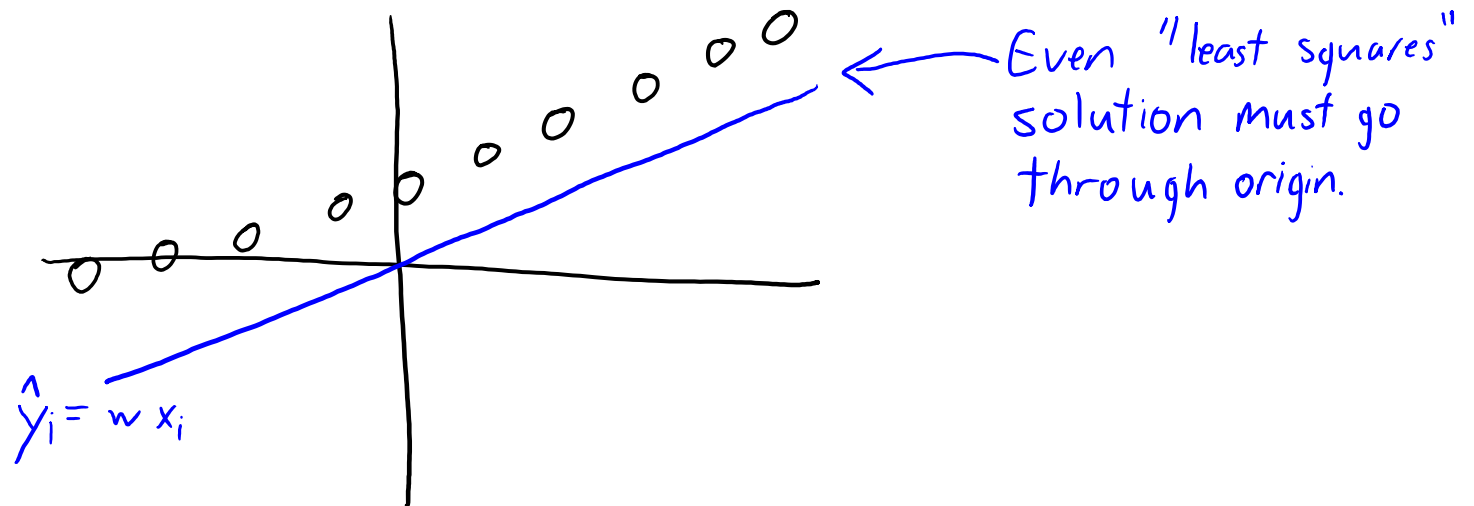                    (Better than KNN and NW at boundaries.)
  - Ensemble methods:
    - Can improve performance by averaging across regression models.

# Adapting Counting/Distance-Based Methods

- 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:
    - http://graphics.cs.cmu.edu/projects/scene-completion
    - Combined with "graph cuts" and "Poisson blending".
  - KNN regression for "voice photoshop":
    - https://www.youtube.com/watch?v=I3l4XLZ59iw
    - 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 $\hat{y}_i = 0$.



Even "least squares" solution must go through origin.

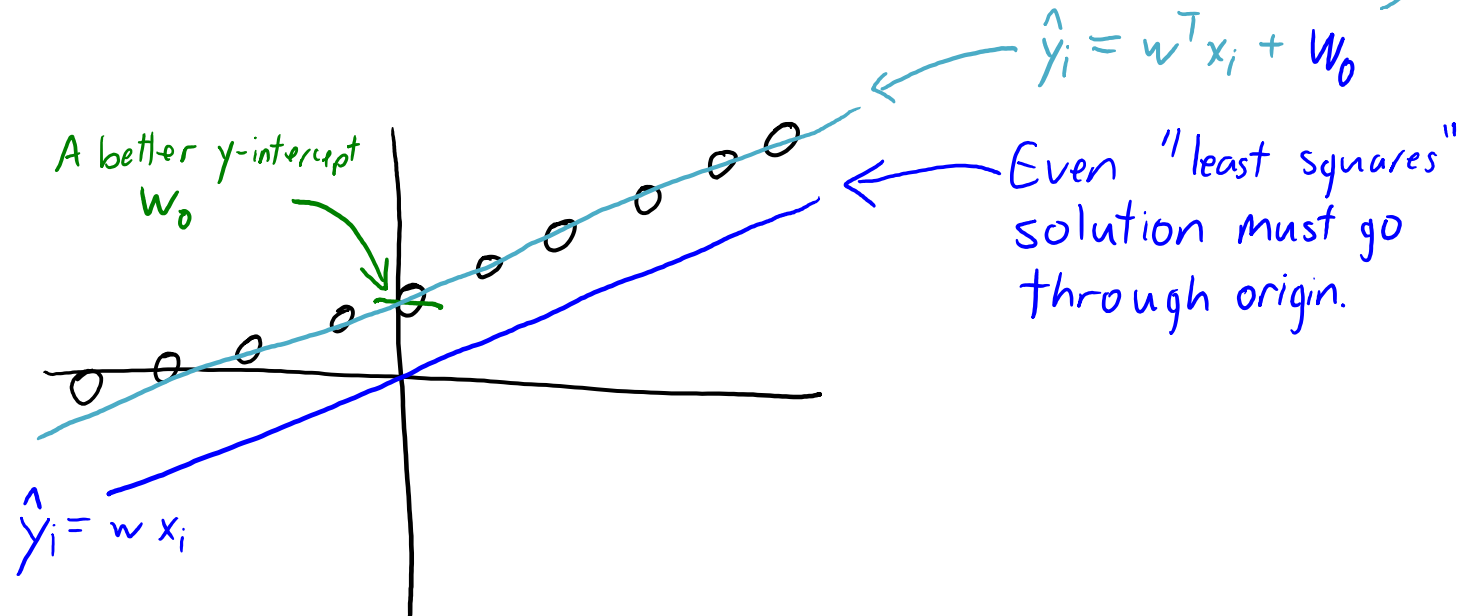$\hat{y}_i = w\,x_i$

# 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
y-intercept
fixes this.

$\hat{y}_i = w^T x_i + w_0$

A better y-intercept
$w_0$

Even "least squares"
solution must go
through origin.

$\hat{y}_i = w\, x_i$

# 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".

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

$$\underbrace{\text{"always 1"}} \quad \underbrace{X}$$

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

$$\hat{y}_i = \underset{w_0}{\underset{\downarrow}{V_1}} \, \underset{1}{\underset{\downarrow}{z_{i1}}} + \underset{w_1}{\underset{\downarrow}{V_2}} \, \underset{x_{i1}}{\underset{\downarrow}{z_{i2}}} = w_0 + w_1 x_{i1}$$

- 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?

$$\hat{y}_i = w_0 + w_1 x_i + w_2 x_i^2$$

- You can do this by changing the features (change of basis):

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

$$\phantom{Z=}\quad y\text{-inf} \quad X \quad x^2$$

- Fit new parameters 'v' under "change of basis": solve $Z^T Z v = Z^T y$.

- It's a linear function of w, but a quadratic function of $x_i$.

$$\hat{y}_i = v^T z_i = v_1 z_{i1} + v_2 z_{i2} + v_3 z_{i3}$$

$$\phantom{\hat{y}_i = v^T z_i =}\quad \underbrace{w_0 \ 1} \qquad \underbrace{w_1 \ x_i} \qquad \underbrace{w_2 \ x_i^2}$$

# Non-Linear Feature Transforms



$\hat{y}_i = w_0 + w_1 x_i$

linear least squares

$y_i = v^T z_i$
$= w_0 + w_1 x_i + w_2 x_i^2$

linear least squares with quadratic basis

To predict on new data $\tilde{X}$, form $\tilde{Z}$ from $\tilde{X}$ and take $y = \tilde{Z}v$

# General Polynomial Features (d=1)

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

$$Z = \begin{bmatrix} 1 & x_1 & (x_1)^2 & \text{-----} & (x_1)^p \\ 1 & x_2 & (x_2)^2 & \text{-----} & (x_2)^p \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & (x_n)^2 & \text{-----} & (x_n)^p \end{bmatrix}$$

- 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

Want 'w' that <u>minimizes</u>

$$f(w) = \frac{1}{2} \sum_{i=1}^{n} (w^T x_i - y_i)^2 = \frac{1}{2} \|Xw - y\|_2^2 = \frac{1}{2}(Xw - y)^T(Xw - y)$$

Let's expand then compute gradient.

$$= \frac{1}{2}\left((Xw)^T - y^T\right)(Xw - y)$$

$$= \frac{1}{2}\left(w^T X^T - y^T\right)(Xw - y)$$

$$= \frac{1}{2}\left(w^T X^T(Xw - y) - y^T(Xw - y)\right)$$

$$= \frac{1}{2}\left(w^T X^T Xw - w^T X^T y - y^T Xw + y^T y\right)$$

$$= \frac{1}{2} w^T X^T Xw - w^T X^T y + \frac{1}{2} y^T y$$

Rule:

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

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

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

$$(A+B)C = AC + BC$$

$$A(B+C) = AB + BC$$

$$\underbrace{a^T A b}_{\text{vector}} = \underbrace{b^T A^T a}_{\text{vector}}$$

Sanity check: all of these are <u>scalars</u>.

# Vector View of Least Squares

- We showed that least squares minimizes:

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

- The ½ and the squaring don't change solution, so equivalent to:

$$f(w) = \| Xw - y \|$$

- 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 _scalars_: $\alpha = 1$, $\beta = 3.5$, $\gamma = \hat{\pi}$

Use _first/last lowercase_ letters for vectors: $w = \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix}$, $x = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$, $y = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$, $a = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$, $b = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$

$\longrightarrow$ Assumed to be _column_-vectors.

Use _first/last uppercase_ letters for matrices: $X, Y, W, A, B$

$\underline{\text{Indices}}$ use $i, j, k$.

$\underline{\text{Sizes}}$ use $m, n, d, p$, and $k$ $\longleftarrow$ hopefully meaning of '$k$' is obvious from context

$\underline{\text{Sets}}$ use $S, T, U, V$

Functions use $f, g$, and $h$.

When I write $x_i$ I 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 model?

# 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", $X^TX$ is invertible, $X^TX$ has non-zero eigenvalues, $\det(X^TX) > 0$.
  - Note that we cannot have independent columns if $d > n$.