# CPSC 340:
# Machine Learning and Data Mining

Nonlinear Regression

# Admin

- <span style="color:red">Midterm</span> on Feb 14 in class.
    - Previous midterms on course homepage.
    - Covers lecture 1-14 (i.e., up to & including today), assignments 1-3.

- Note to self: we need to finite the previous lecture before starting.

# Summary of Last Lecture

1. Error functions:
   - Squared error is sensitive to outliers.
   - Absolute ($L_1$) error and Huber error are more robust to outliers.
   - Brittle ($L_\infty$) error is more sensitive to outliers.
2. $L_1$ and $L_\infty$ error functions are convex but non-differentiable:
   - Finding 'w' that minimizes these errors is harder than squared error.
3. We can approximate these with convex differentiable functions:
   - $L_1$ can be approximated with Huber.
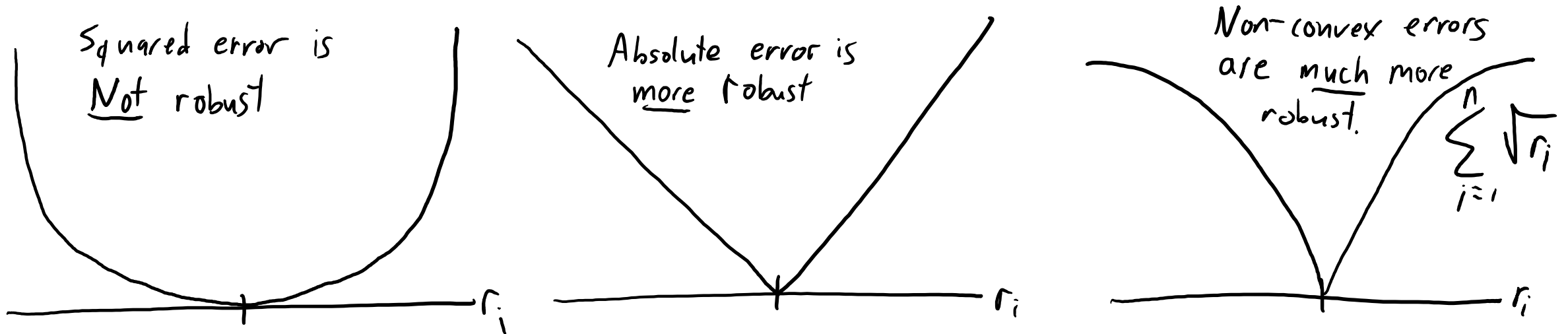   - $L_\infty$ can be approximated with log-sum-exp.
4. Gradient descent finds stationary point of differentiable function.
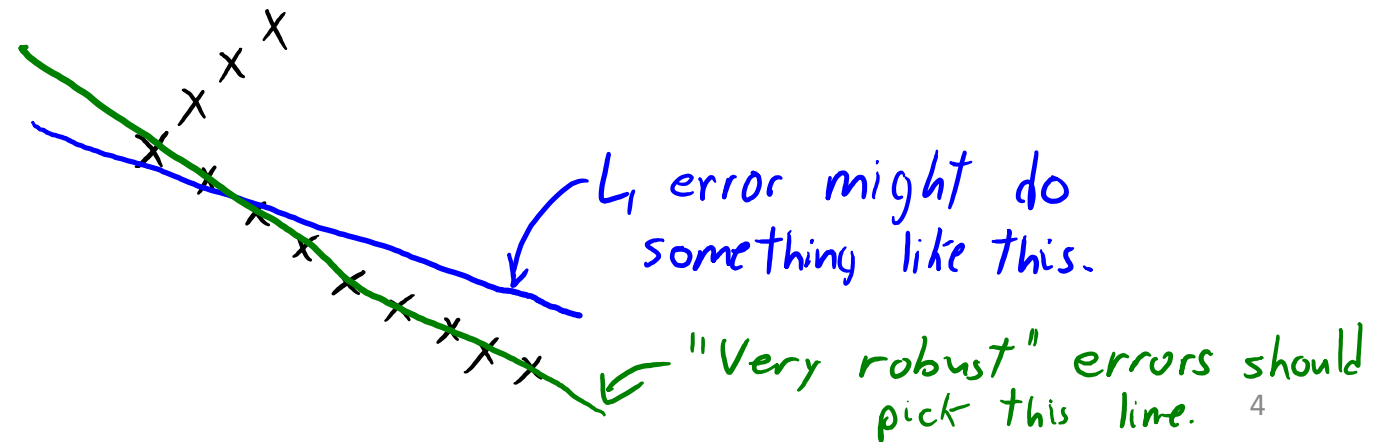   - "Stationary point" == "critical point" == "a value of 'w' where $\nabla f(w) = 0$".
5. For convex functions, any stationary point is a global minimum.
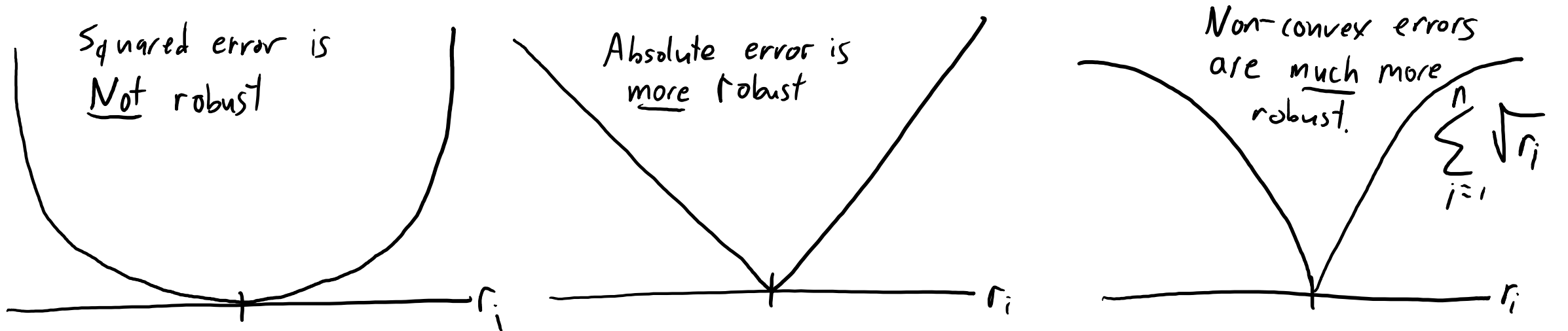   - So gradient descent finds global minimum.

# Very Robust Regression

Squared error is __Not__ robust

Absolute error is __more__ robust

Non-convex errors are __much__ more __robust__.

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

$r_i$    $r_i$    $r_i$

- Non-convex errors can be very robust:
    - Not influenced by outlier groups.

$L_1$ error might do something like this.

"Very robust" errors should pick this line.

# Very Robust Regression

Squared error is __Not__ robust

Absolute error is __more__ robust

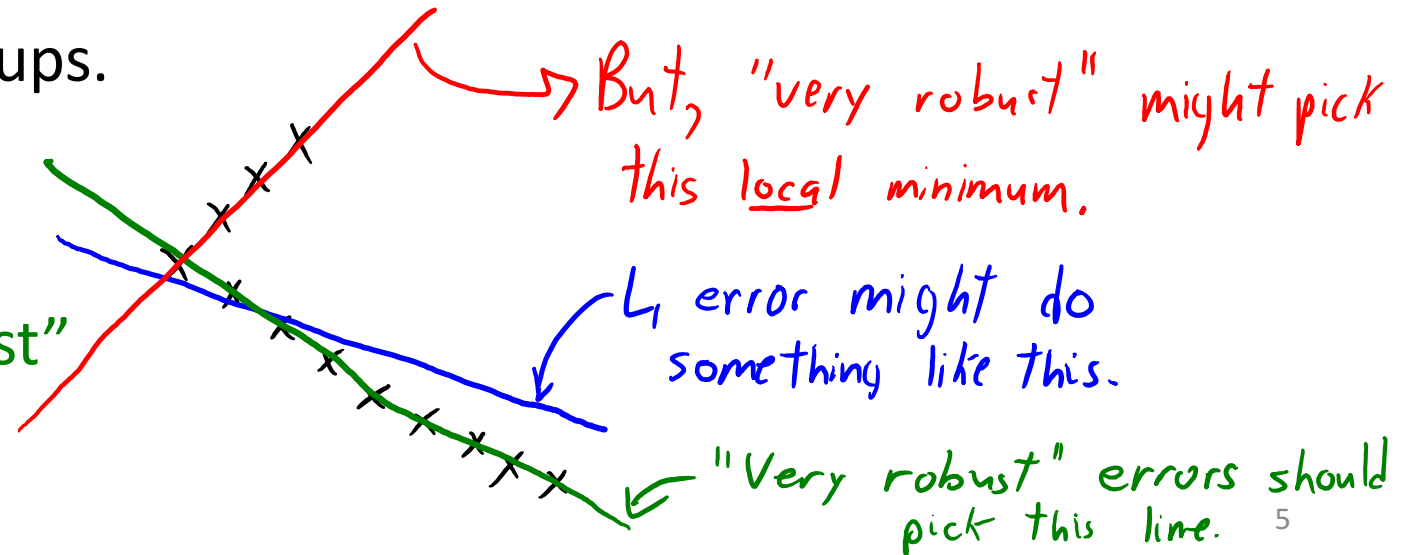Non-convex errors are __much__ more __robust__.

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

$r_i$

$r_i$

$r_i$

- Non-convex errors can be very robust:
  - Not influenced by outlier groups.
  - But non-convex, so finding global minimum is hard.
  - Absolute value is "most robust" convex loss function.

But, "very robust" might pick this __local__ minimum.

$L_1$ error might do something like this.

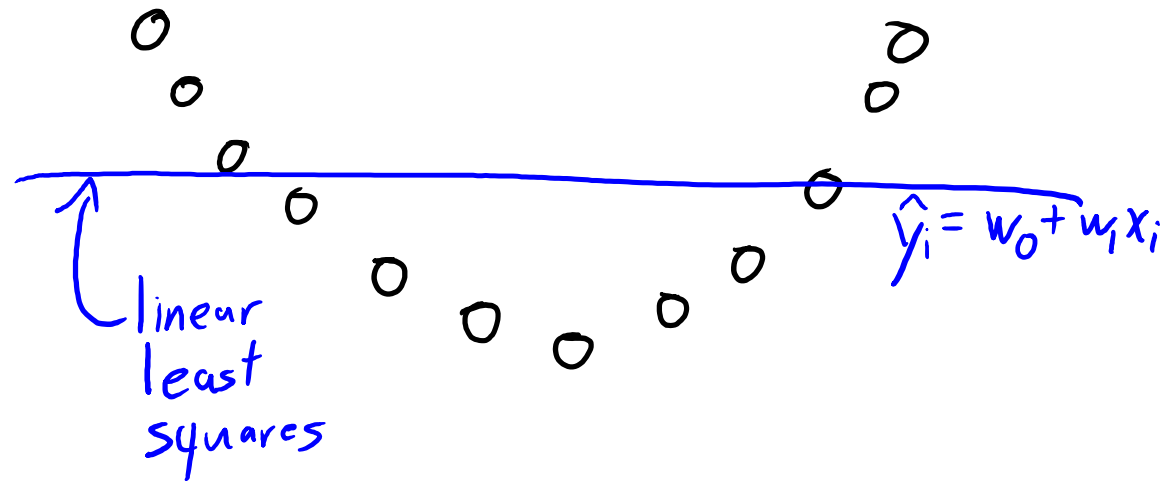"Very robust" errors should pick this line.

5

(pause)

# Nonlinear regression

- We can adapt classification methods to perform regression.
- E.g. decision tree regression, KNN regression.
- See bonus slides for more details.
- We will focus on direct extensions of linear regression.

# Motivation: Limitations of Linear Models

- On many datasets, y$_i$ is not a linear function of x$_i$.



$$\hat{y_i} = w_0 + w_1 x_i$$

linear
least
squares

- 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?

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

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

$$X = \begin{bmatrix} 0.2 \\ -0.5 \\ 1 \\ 4 \end{bmatrix} \qquad Z = \begin{bmatrix} 1 & 0.2 & (0.2)^2 \\ 1 & -0.5 & (-0.5)^2 \\ 1 & 1 & (1)^2 \\ 1 & 4 & (4)^2 \end{bmatrix}$$
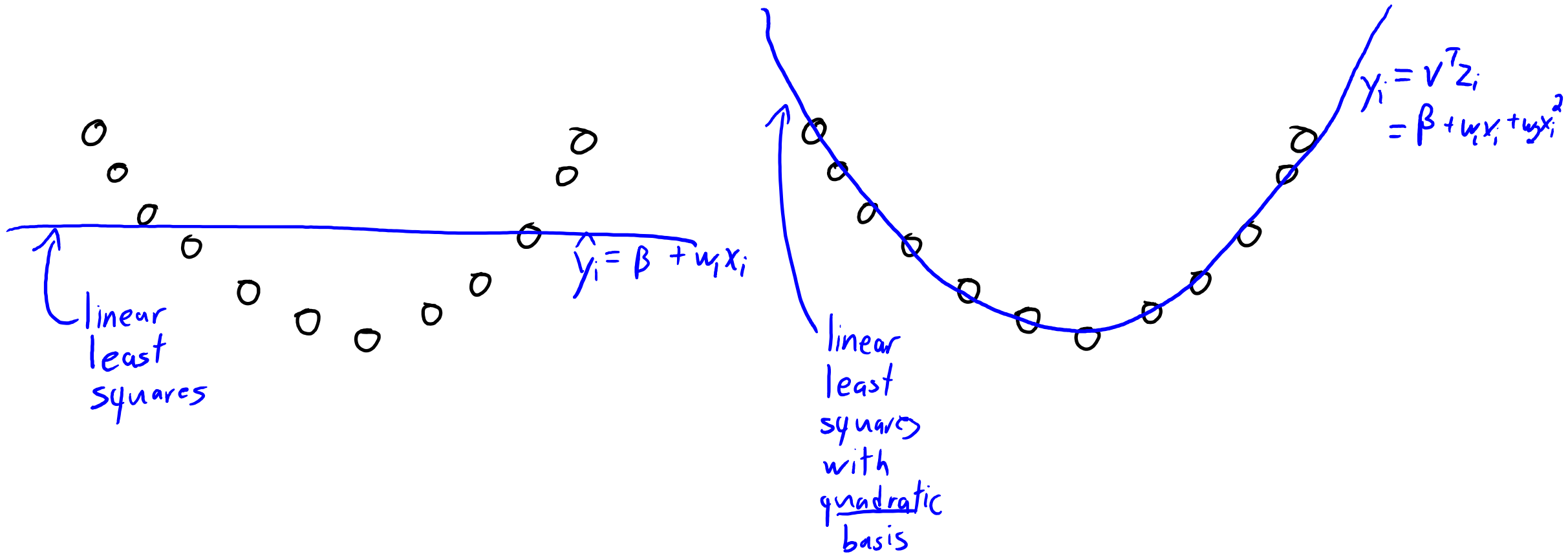
$$\begin{aligned} \hat{y}_i &= v^T z_i \\ &= v_1 z_{i0} + v_2 z_{i1} + v_3 z_{i2} \\ &\simeq \beta + w_1 x_i + w_2 x_i^2 \end{aligned}$$

y-inf    X    $x^2$

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

- Fit using normal equations with Z instead of X: $v = (Z^T Z)^{-1}(Z^T y)$

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

# Non-Linear Feature Transforms



linear
least
squares

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

linear
least
squares
with
quadratic
basis

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

# General Polynomial Features (d=1)

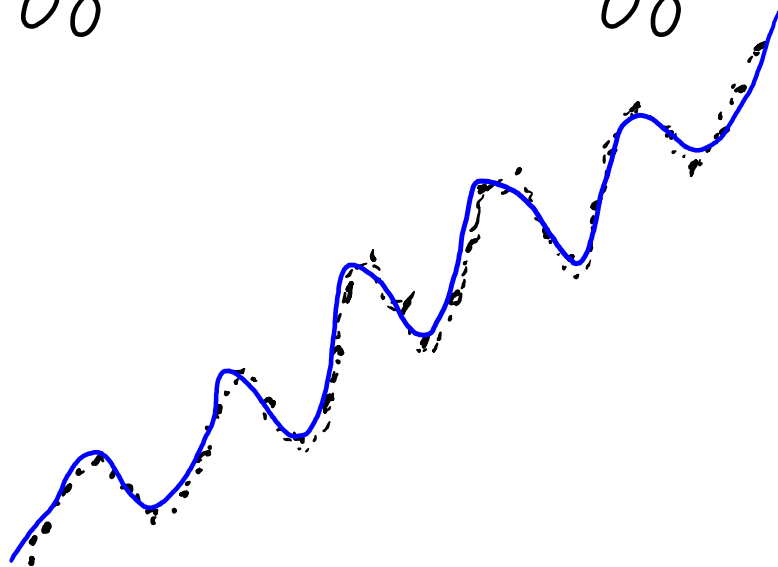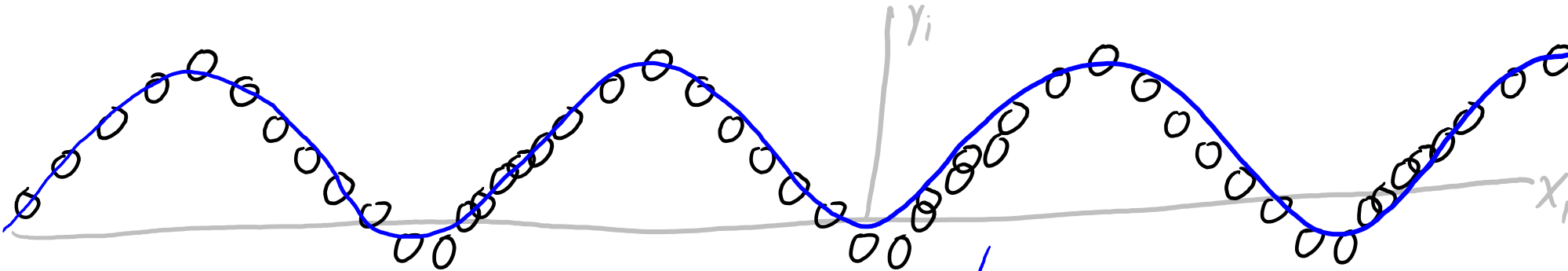- We can have a polynomial of degree 'p' by using these features:

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

- There are polynomial basis functions that are numerically nicer:
  - E.g., Lagrange polynomials (see CPSC 303).
- If you have more than one feature, you include interactions:
  - With p=2, in addition to $(x_{i1})^2$ and $(x_{i2})^2$ you would include $x_{i1}x_{i2}$.

# Jupyter notebook demo

# Beyond Polynomial Transformations

- Polynomials are not the only possible transformation:
  - Exponentials, logarithms, trigonometric functions, etc.
  - The right non-linear transform will vastly improve performance.

For periodic data
we might use

$$Z = \begin{bmatrix} \sin(x_1) \\ \sin(x_2) \\ \vdots \\ \sin(x_n) \end{bmatrix}$$

You can have different types
of bases

$$Z = \begin{bmatrix} x_1 & \sin(6x_1) \\ x_2 & \sin(6x_2) \\ \vdots & \vdots \\ x_n & \sin(6x_n) \end{bmatrix}$$
$\underbrace{\phantom{x_n}}_{\text{linear}} \underbrace{\phantom{\sin(6x_n)}}_{\text{periodic}}$

$\hat{y}_i = v^T z_i$

$\quad = w_1 \sin(x_i)$

# Parametric vs. Non-Parametric Transforms

- We've been using linear models with <span style="color:blue">polynomial bases</span>:

$$y_i = w_0 \,\square\, + w_1 \,\square\, + w_2 \,\square\, + w_3 \,\square\, + w_4 \,\square$$

- But polynomials are not the only <span style="color:green">possible bases</span>:
  - Exponentials, logarithms, trigonometric functions, etc.
  - The <span style="color:green">right basis will vastly improve performance</span>.
  - If we use the wrong basis, our accuracy is limited even with lots of data.
  - But the <span style="color:red">right basis may not be obvious</span>.

# Parametric vs. Non-Parametric Transforms

- We've been using linear models with polynomial bases:

$$y_i = w_0 \; \boxed{\phantom{xxx}} + w_1 \; \boxed{\phantom{xxx}} + w_2 \; \boxed{\phantom{xxx}} + w_3 \; \boxed{\phantom{xxx}} + w_4 \; \boxed{\phantom{xxx}}$$

- Alternative is non-parametric bases:
  - Size of basis (number of features) grows with 'n'.
  - Model gets more complicated as you get more data.
  - Can model complicated functions where you don't know the right basis.
    - With enough data.
  - Classic example is "Gaussian RBFs".

# Gaussian RBFs: A Sum of "bumps"



$$y_i = w_0 \boxed{\quad} + w_1 \boxed{\quad} + w_2 \boxed{\quad} + w_3 \boxed{\quad} + w_4 \boxed{\quad}$$

Polynomial basis represents function as sum of g**lobal** polynomials.

$$y_i = w_0 \boxed{\quad} + w_1 \boxed{\quad} + w_2 \boxed{\quad} + w_3 \boxed{\quad} + w_4 \boxed{\quad}$$

Gaussian RBFs represent function as sum of **local** "bumps"

- Gaussian RBFs are universal approximators (compact subets of $\mathbb{R}^d$)
  - Enough bumps can approximate any continuous function to arbitrary precision.
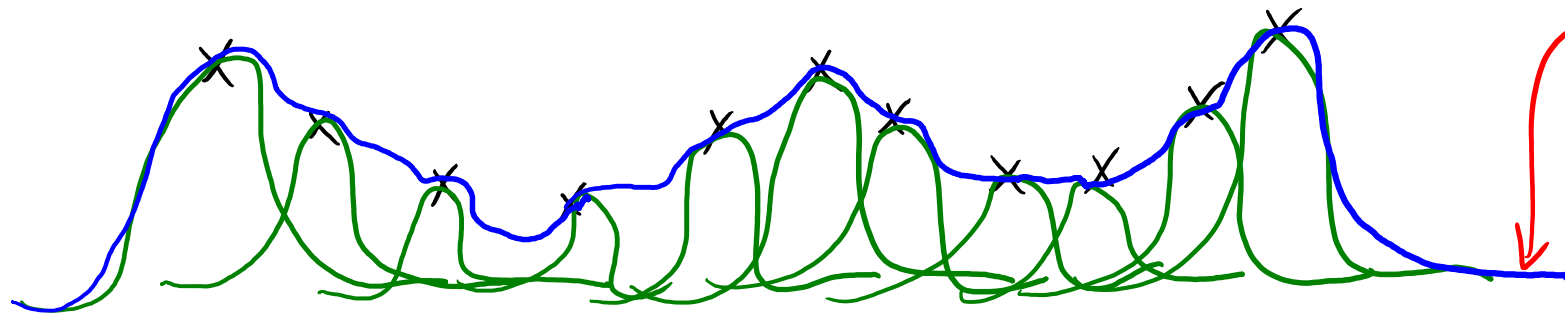  - Achieve optimal test error as 'n' goes to infinity.

# Gaussian RBFs: A Sum of "Bumps"

- Polynomial fit:



polynomial basis becomes polynomial away from data

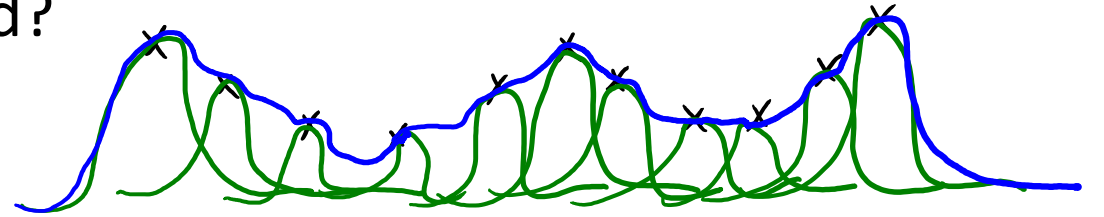- Constructing a function from bumps:



Gaussian RBFs go to zero away from data.

- Bonus slides: challenges of "far from data" (and future) predictions.
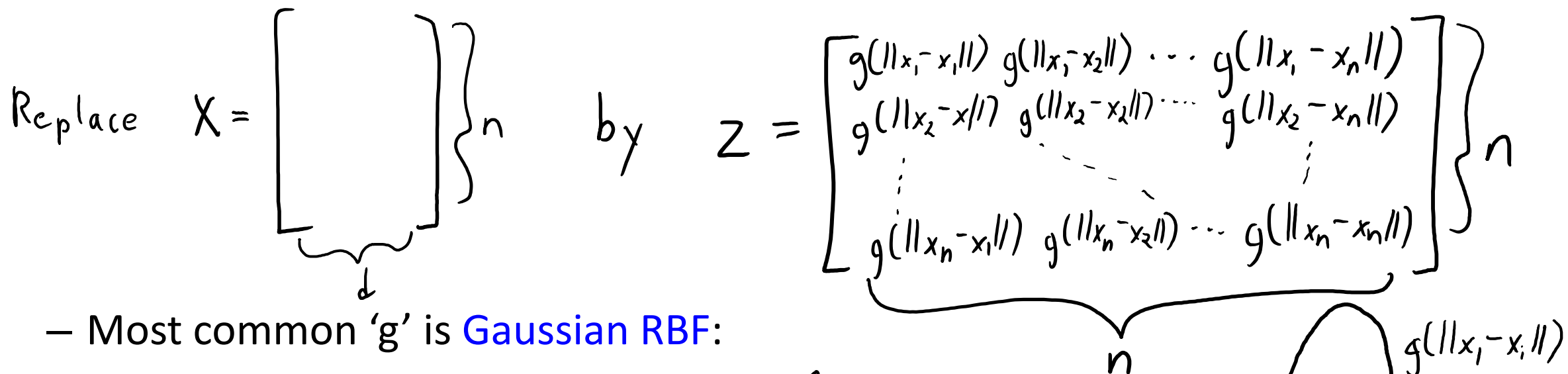
# Gaussian RBF Parameters

- Some obvious questions:

  1. How many bumps should we use?

  2. Where should the bumps be centered?

  3. How high should the bumps go?

  4. How wide should the bumps be?



- The usual answers:

  1. We use 'n' bumps (non-parametric basis).

  2. Each bump is centered on one training example $x_i$.

  3. Fitting regression weights 'w' gives us the heights (and signs).

  4. The width is a hyper-parameter (narrow bumps == complicated model).
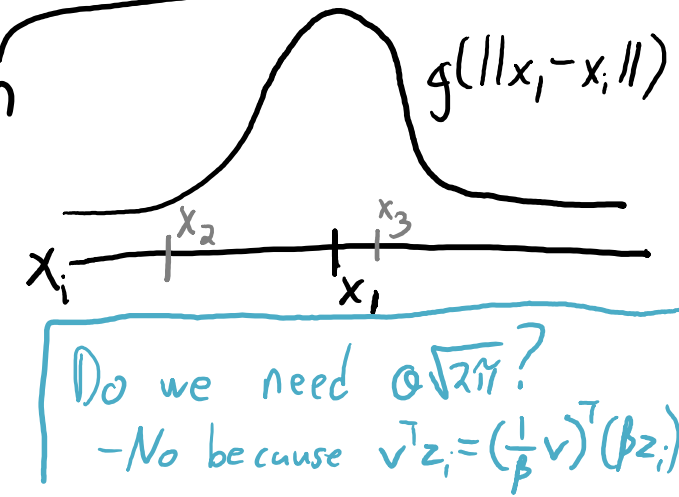
# Gaussian RBFs: Formal Details

- What is a radial basis functions (RBFs)?
  - A set of non-parametric bases that depend on distances to training points.

$$Replace \quad X = \begin{bmatrix} \phantom{xxx} \end{bmatrix} \Big\}n \quad by \quad Z = \begin{bmatrix} g(\|x_1-x_1\|) & g(\|x_1-x_2\|) & \cdots & g(\|x_1-x_n\|) \\ g(\|x_2-x_1\|) & g(\|x_2-x_2\|) & \cdots & g(\|x_2-x_n\|) \\ \vdots & & & \vdots \\ g(\|x_n-x_1\|) & g(\|x_n-x_2\|) & \cdots & g(\|x_n-x_n\|) \end{bmatrix} \Big\}n$$

$d$

$n$

$g(\|x_1-x_i\|)$

  - Most common 'g' is Gaussian RBF:

$$g(\varepsilon) = exp\left(-\frac{\varepsilon^2}{2\sigma^2}\right)$$

  - Variance $\sigma^2$ is a hyper-parameter controlling "width".
    - This affects fundamental trade-off (set it using a validation set).

$x_i$ $x_2$ $x_3$ $x_1$

Do we need $\sigma\sqrt{2\pi}$?
- No because $v^T z_i = (\frac{1}{\beta}v)^T(\beta z_i)$

# Gaussian RBFs: Formal Details

- What is a radial basis functions (RBFs)?
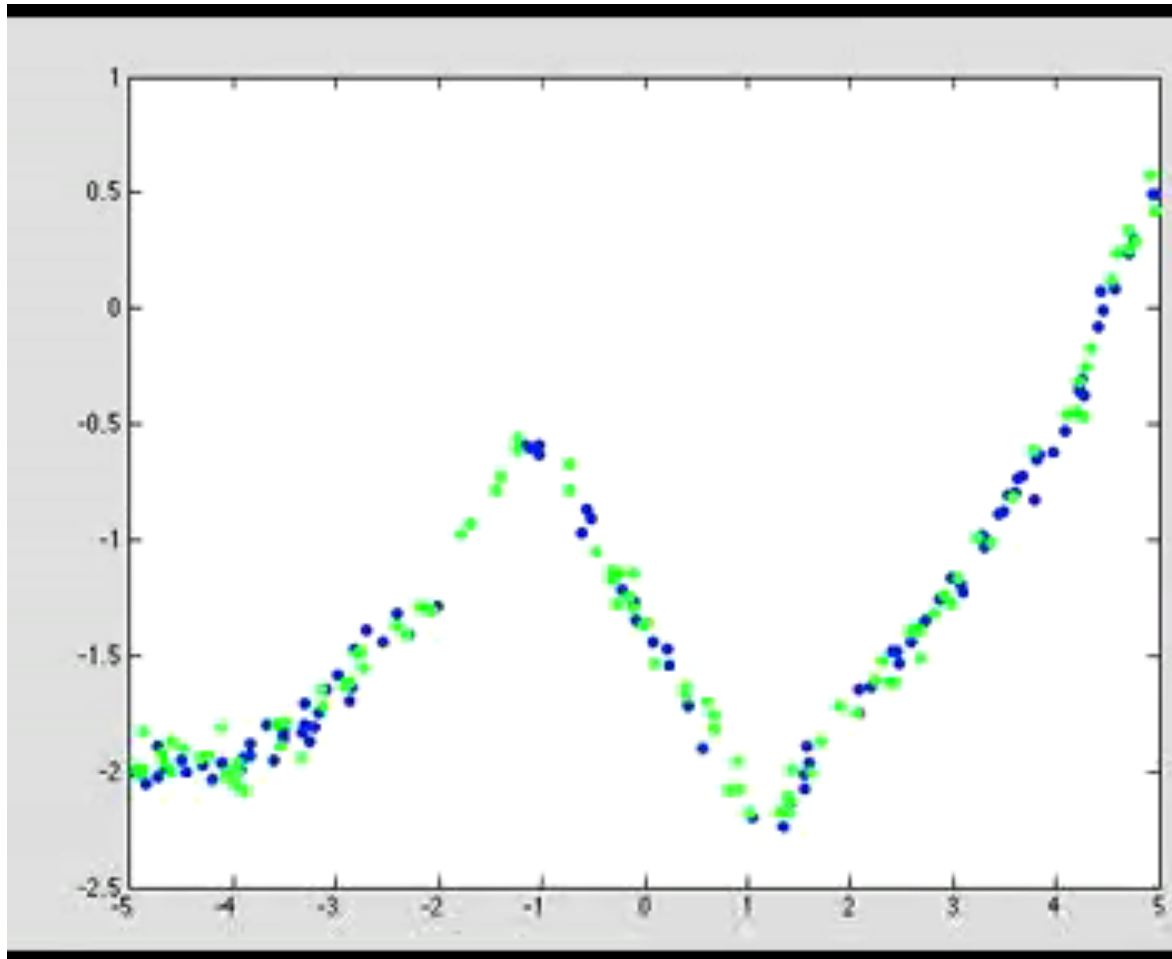  - A set of non-parametric bases that depend on distances to training points.

$$\text{Replace} \quad X = \begin{bmatrix} \quad \\ \quad \\ \quad \end{bmatrix} \Big\} n \quad \text{by} \quad Z = \begin{bmatrix} g(\|x_1 - x_1\|) & g(\|x_1 - x_2\|) & \cdots & g(\|x_1 - x_n\|) \\ g(\|x_2 - x_1\|) & g(\|x_2 - x_2\|) & \cdots & g(\|x_2 - x_n\|) \\ \vdots & & & \\ g(\|x_n - x_1\|) & g(\|x_n - x_2\|) & \cdots & g(\|x_n - x_n\|) \end{bmatrix} \Big\} n$$

$$\underbrace{\qquad}_{d} \qquad\qquad \underbrace{\qquad\qquad\qquad\qquad}_{n}$$

$$\text{To make predictions on} \quad \tilde{X} = \begin{bmatrix} \quad \\ \quad \\ \quad \end{bmatrix} \Big\} t \quad \text{use} \quad \tilde{Z} = \begin{bmatrix} \\ g(\|\tilde{x}_i - x_j\|) \\ \\ \end{bmatrix} \Big\} t$$

$$\underbrace{\qquad}_{d} \qquad\qquad\qquad \underbrace{\qquad\qquad\qquad}_{n}$$

Number of "features" is number of training examples

# Non-Parametric Basis: RBFs

- Least squares with Gaussian RBFs for different σ values:



Could add __bias__ and linear basis:

$$Z = \begin{bmatrix} 1 & - x_1 - & g(\|x_1 - x_1\|) & \cdots & g(\|x_1 - x_n\|) \\ 1 & - x_2 - & & & \vdots \\ 1 & - x_3 - & \vdots & & \\ \vdots & \vdots & & & \vdots \\ 1 & - x_n - & g(\|x_1 - x_n\|) & \cdots & g(\|x_n - x_n\|) \end{bmatrix}$$

$\underbrace{\quad}_{1} \underbrace{\quad}_{d} \underbrace{\qquad\qquad}_{n}$

This reverts to linear regression instead of 0 away from data.

# Summary

- Tree/probabilistic/non-parametric/ensemble regression methods.
- Non-linear transforms:
  - Allow us to model non-linear relationships with linear models.
  - Polynomial features are a parametric example
  - RBF features are a non-parametric example

# Complexity Penalties

- The next 3 slides are a previous of next week's topics.

# Finding the "True" Model

- What if our goal is find the "true" model?
  - We believe that $y_i$ really is a polynomial function of $x_i$.
  - We want to find the degree of the polynomial 'p'.

- Should we choose the 'p' with the lowest training error?
  - No, this will pick a 'p' that is way too large.
    (training error always decreases as you increase 'p')

# Finding the "True" Model

- What if our goal is find the "true" model?
  - We believe that $y_i$ really is a polynomial function of $x_i$.
  - We want to find the degree of the polynomial 'p'.

- Should we choose the 'p' with the lowest validation error?
  - This will also often choose a 'p' that is too large.

  - Even if true model has p=2, this is a special case of a degree-3 polynomial.
  - If 'p' is too big then we overfit, but might still get a lower validation error.
    - Another example of optimization bias.

# Complexity Penalties

- There are a lot of "scores" people use to find the "true" model.

- Basic idea behind them: put a penalty on the model complexity.
  - Want to **fit the data and have a simple model**.

- For example, minimize training error plus the degree of polynomial.

$$\text{Let } Z_p = \begin{bmatrix} 1 & x_1 & (x_1)^2 & \cdots & (x_1)^p \\ 1 & x_2 & (x_2)^2 & \cdots & (x_2)^p \\ 1 & x_3 & (x_3)^2 & \cdots & (x_3)^p \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & (x_n)^2 & \cdots & (x_n)^p \end{bmatrix}$$

Find 'p' that minimizes:

$$\text{score}(p) = \frac{1}{2}\|Z_p v - y\|^2 + p$$
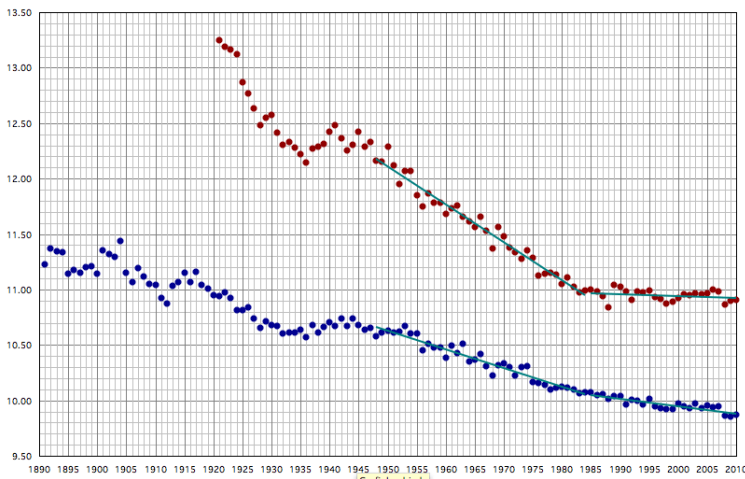
train error for best 'v' with this basis.

degree of polynomial

  - If we use p=4, use "training error plus 4" as error.

- If two 'p' values have similar error, this prefers the smaller 'p'.
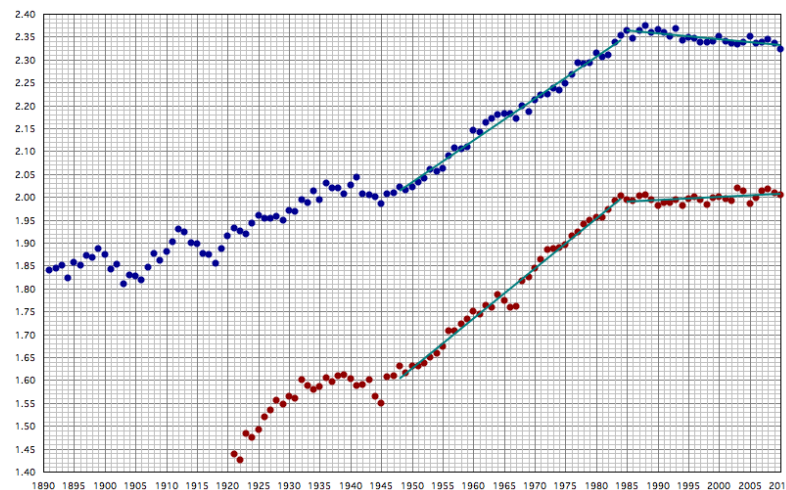
# Motivation: Non-Linear Progressions in Athletics

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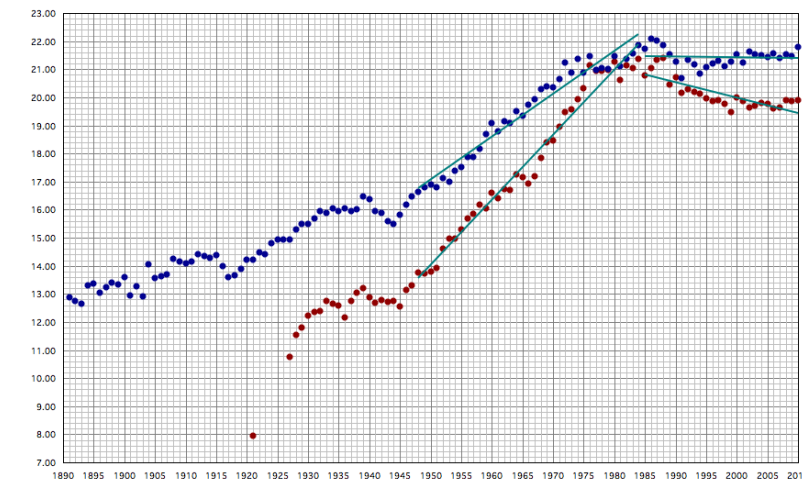


100m PROGRESSION MEN AND WOMEN (mean of top ten)



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



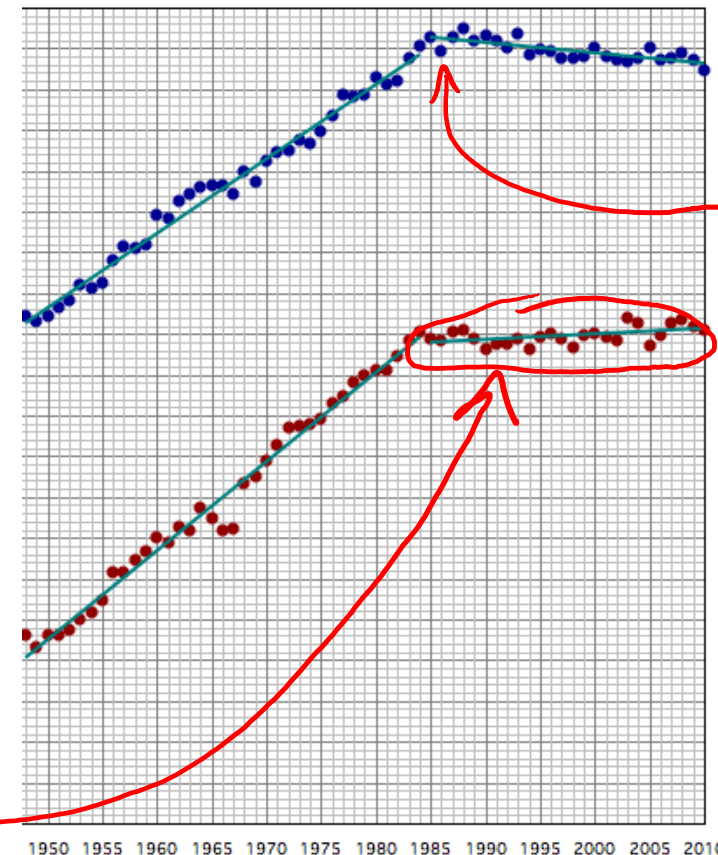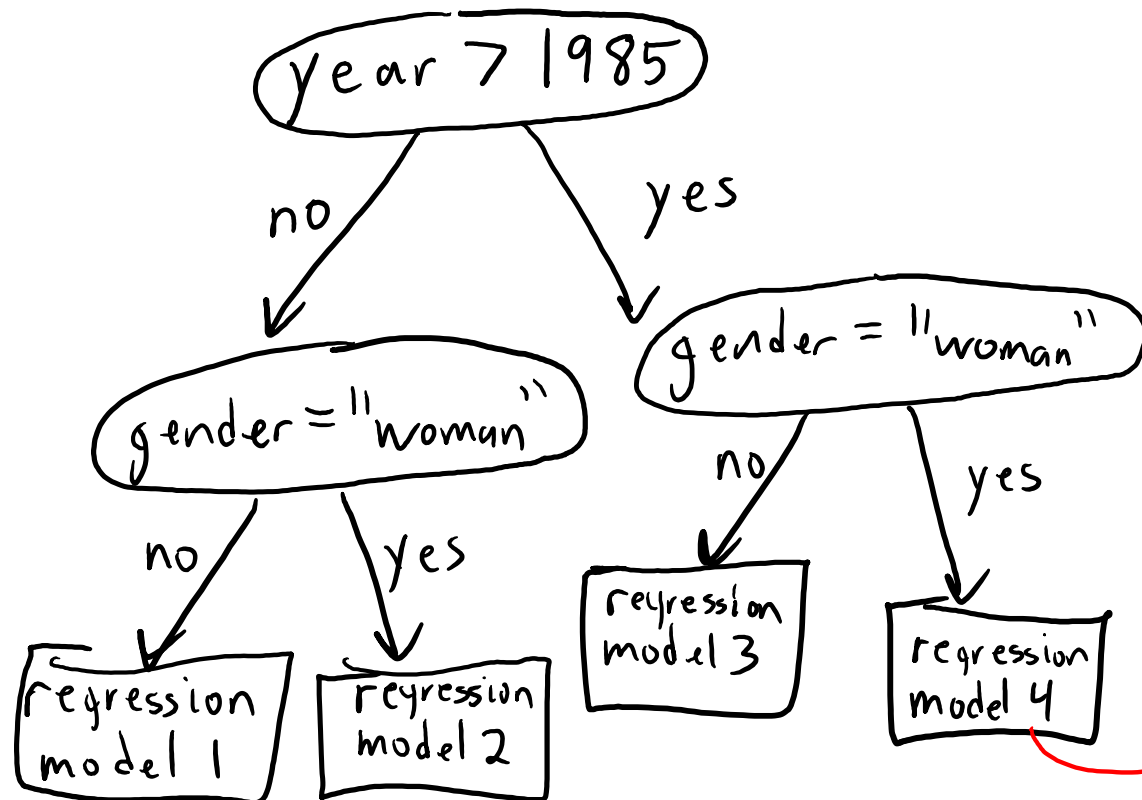SHOT PUT PROGRESSION MEN (7.26 kg) AND WOMEN (4 kg) (mean of top ten)

27

# Adapting Counting/Distance-Based Methods

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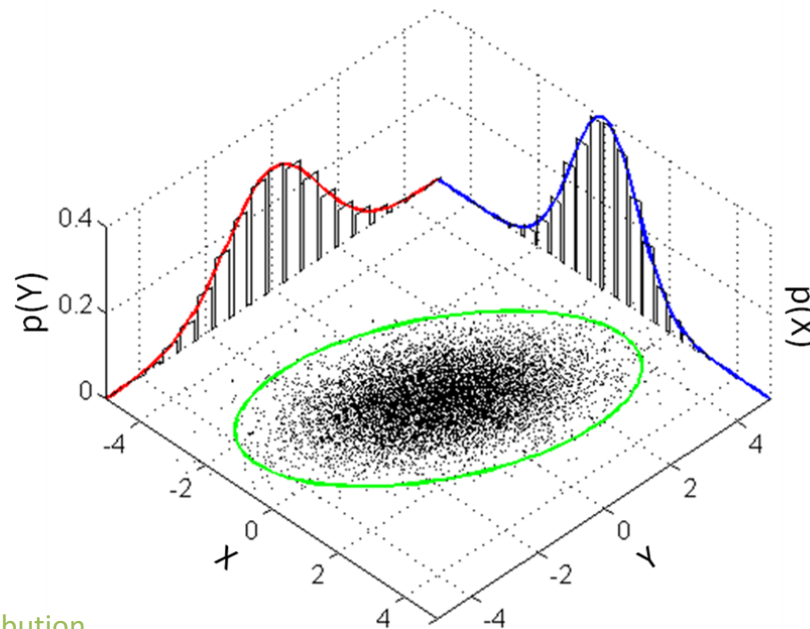
# Adapting Counting/Distance-Based Methods

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

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

# 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.
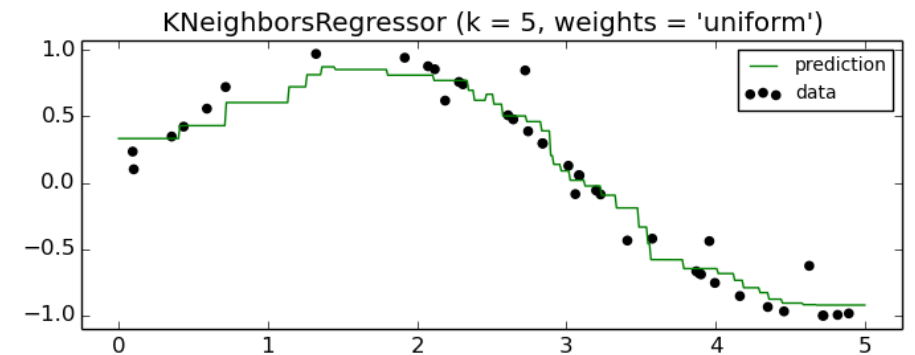    - CPSC 540.

30

# 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 $x_i$.
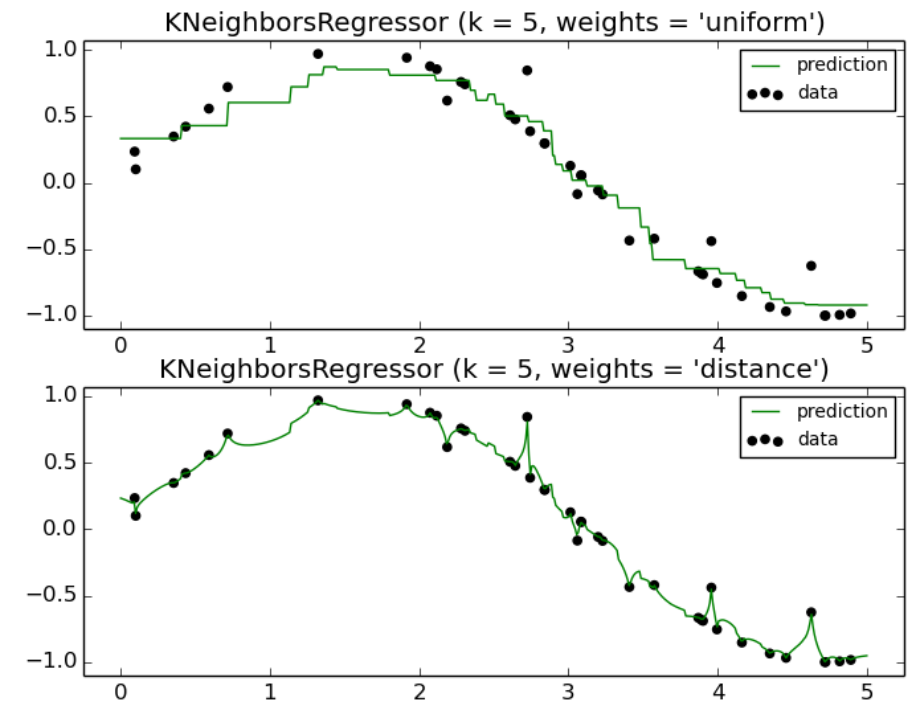      - Return the mean of the corresponding $y_i$.

# 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" $w_{ij}$.



http://scikit-learn.org/stable/modules/neighbors.html
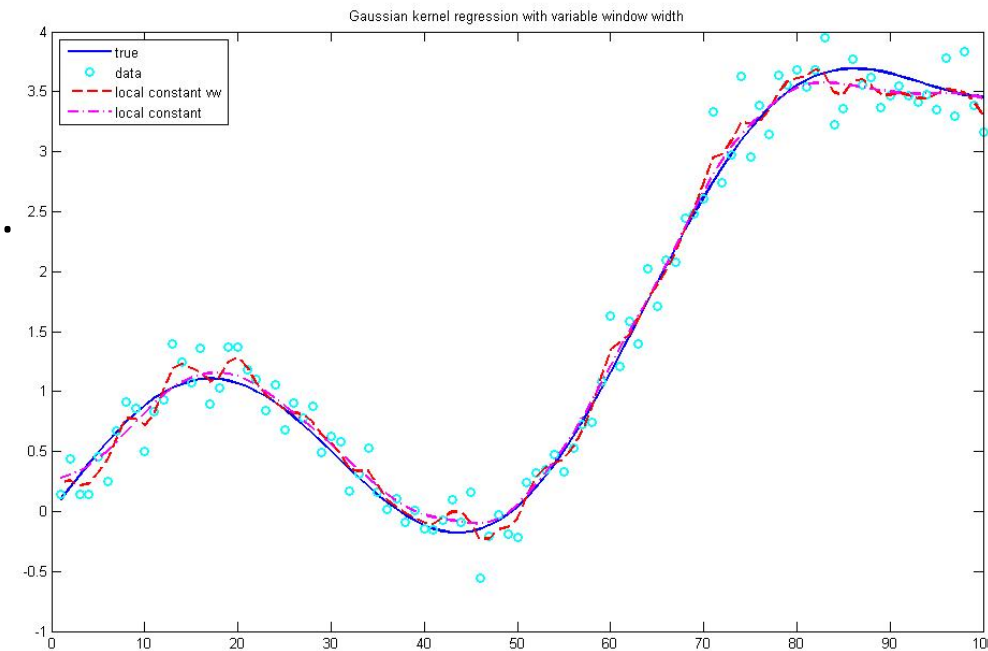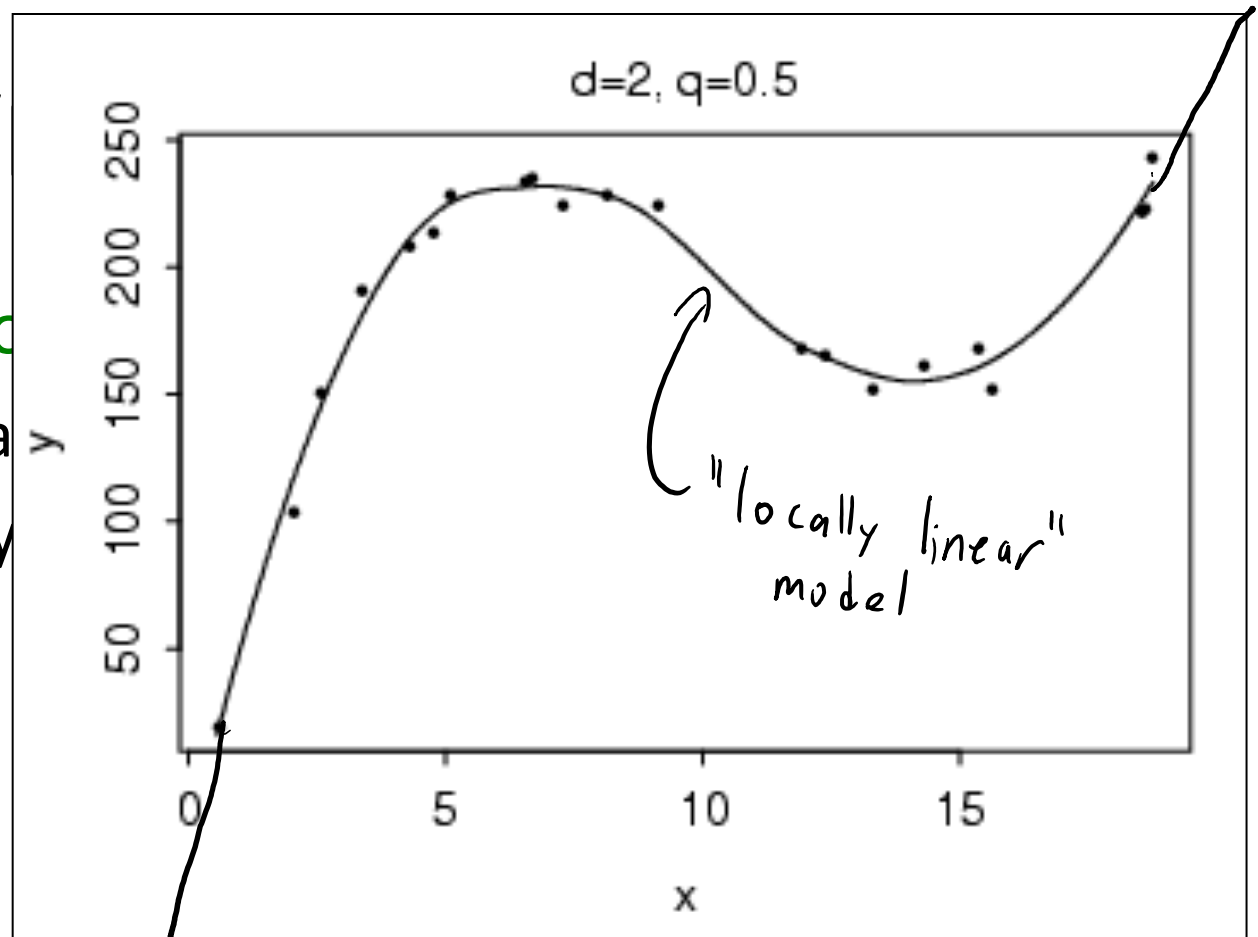
# Adapting Counting/Distance-Based Methods

- 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

# Adapting Counting/

- We can adapt our classificatio
  - 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.

http://www.itl.nist.gov/div898/handbook/pmd/section4/pmd423.htm