# COMS 4771
# Regression

Nakul Verma

# Last time…

- Support Vector Machines

- Maximum Margin formulation

- Constrained Optimization

- Lagrange Duality Theory

- Convex Optimization

- SVM dual and Interpretation

- How get the optimal solution

# Learning more Sophisticated Outputs

So far we have focused on classification $f : X \rightarrow \{1, ..., k\}$
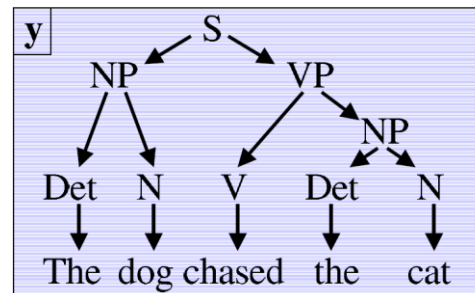
What about **other outputs**?

- PM$_{2.5}$ (pollutant) particulate matter exposure estimate:

   **Input:** # cars, temperature, etc.          **Output:** 50 ppb
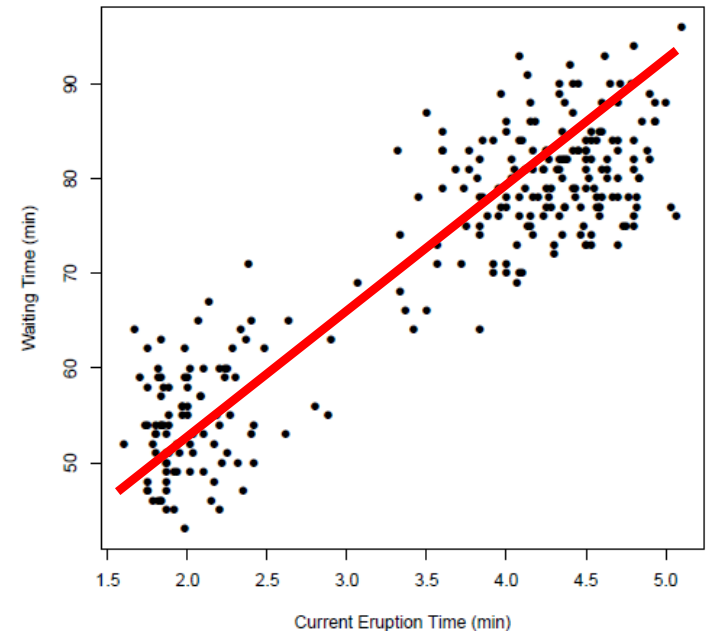
- Pose estimation

- Sentence structure estimate:

# Regression

We'll focus on problems with real number outputs (regression problem):

$$f : X \to \mathbf{R}$$

Example:

Next eruption time of old faithful geyser (at Yellowstone)

# Regression Formulation for the Example

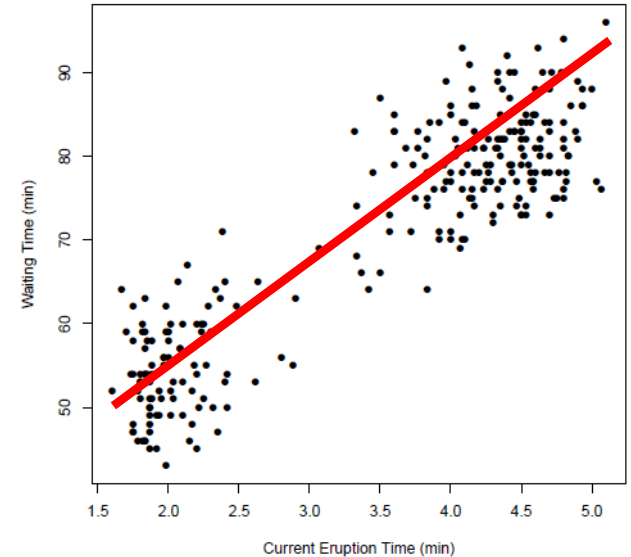Given *x,* want to predict an estimate *ŷ* of *y*, which minizes the discrepancy (*L*) between *ŷ* and *y*.

$$L(\hat{y}; y) := |\hat{y} - y| \qquad \textit{Absolute error}$$

$$:= (\hat{y} - y)^2 \qquad \textit{Squared error}$$

Loss



A **linear predictor** *f*, can be defined by the slope *w* and the intercept $w_0$ :

$$\hat{f}(\vec{x}) := \vec{w} \cdot \vec{x} + w_0$$

which minimizes the prediction loss.

*How is this different from **classification**?*

$$\min_{w, w_0} \mathbb{E}_{\vec{x}, y} \left[ L(\hat{f}(\vec{x}), y) \right]$$

# Parametric vs non-parametric Regression

If we assume a particular form of the regressor:
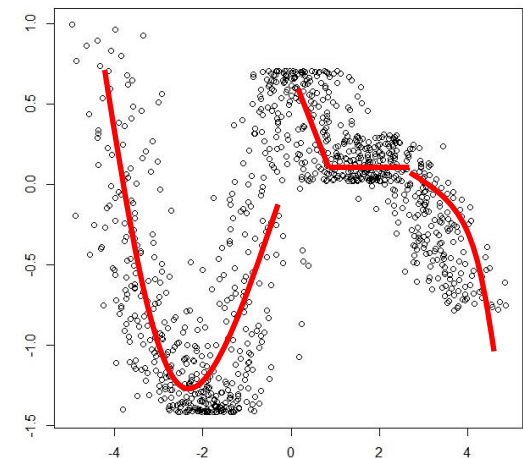
*Parametric regression*

*Goal: to learn the parameters which yield the minimum error/loss*



If no specific form of regressor is assumed:

*Non-parametric regression*

*Goal: to learn the predictor directly from the input data that yields the minimum error/loss*

# Linear Regression

Want to find a **linear predictor** *f*, i.e., *w* (intercept $w_0$ absorbed via lifting):

$$\hat{f}(\vec{x}) := \vec{w} \cdot \vec{x}$$

which minimizes the prediction loss over the population.

$$\min_{\vec{w}} \mathbb{E}_{\vec{x},y} \left[ L(\hat{f}(\vec{x}), y) \right]$$

*(Geometrically)*

We estimate the parameters by minimizing the corresponding loss on the training data:

$$\arg\min_{w} \frac{1}{n} \sum_{i=1}^{n} \left[ L(\vec{w} \cdot \vec{x}_i, y_i) \right]$$

$$= \arg\min_{w} \frac{1}{n} \sum_{i=1}^{n} \left( \vec{w} \cdot \vec{x}_i - y_i \right)^2$$

*for squared error*

# Linear Regression: Learning the Parameters

Linear predictor with squared loss:

$$\arg\min_{w} \frac{1}{n} \sum_{i=1}^{n} \left( \vec{w} \cdot \vec{x}_i - y_i \right)^2$$

$$= \arg\min_{w} \left\| \begin{pmatrix} \dots x_1 \dots \\ \dots x_i \dots \\ \dots x_n \dots \end{pmatrix} \begin{pmatrix} w \end{pmatrix} - \begin{pmatrix} y_1 \\ y_i \\ y_n \end{pmatrix} \right\|^2$$

$$= \arg\min_{w} \left\| X\vec{w} - \vec{y} \right\|_2^2$$

*Unconstrained problem!*

*Can take the gradient and examine the stationary points!*

*Why need not check the second order conditions?*

# Linear Regression: Learning the Parameters

Thus best fitting *w*:
$$\frac{\partial}{\partial \vec{w}} \left\| X\vec{w} - \vec{y} \right\|^2 = 2X^{\mathsf{T}}(X\vec{w} - \vec{y})$$

*At a stationarity*
$$X^{\mathsf{T}}X\vec{w} = X^{\mathsf{T}}\vec{y}$$
***This system is always consistent!***

*why?*
$$\vec{y} = \vec{y}_{\mathrm{col}(X)} + \vec{y}_{\mathrm{null}(X^{\mathsf{T}})}$$
(via orth decomp of y)

*thus*
$$X^{\mathsf{T}}\vec{y} = X^{\mathsf{T}}\vec{y}_{\mathrm{col}(X)}$$

*since*
$$\vec{y}_{\mathrm{col}(X)} = \sum_i w_i \ddot{x}_i$$
(for some coefficients $w_i$, where $\ddot{x}_i$ are columns of X)

*define*
$$\vec{w} := \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$$

*Then*
$$X^{\mathsf{T}}X\vec{w} = X^{\mathsf{T}}(X\vec{w}) = X^{\mathsf{T}}\left(\sum_i w_i \ddot{x}_i\right)$$
$$= X^{\mathsf{T}}\vec{y}_{\mathrm{col}(X)} = X^{\mathsf{T}}\vec{y}$$

*How can we find this w*
*(ie the coefficients $w_i$)*
*which satisfies the system?*

$$\vec{w}_{\mathrm{ols}} = (X^{\mathsf{T}}X)^{\dagger}X^{\mathsf{T}}\vec{y}$$

Pseudo-inverse

***Also called the Ordinary***
***Least Squares (OLS)***

*The solution is unique and*
*stable when $X^TX$ is invertible*

*What is the interpretation of this solution?*

# Linear Regression: Geometric Viewpoint

Consider the **column space** view of data **X**:

$$\begin{pmatrix} \dots x_1 \dots \\ \dots x_i \dots \\ \dots x_n \dots \end{pmatrix} \qquad \ddot{x}_1, \dots, \ddot{x}_d \in \mathbf{R}^n$$

Find a *w*, such that the linear combination of **minimizes**

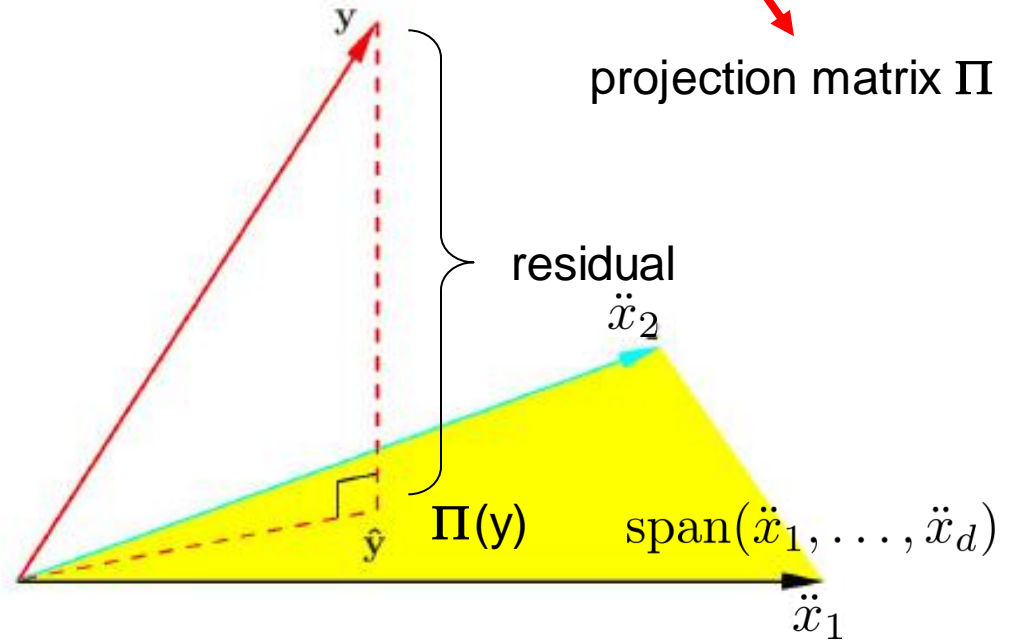$$\frac{1}{n} \left\| \vec{y} - \sum_{i=1}^{d} w_i \ddot{x}_i \right\|^2 =: \text{residual}$$

$$\hat{y} = X \vec{w}_{\text{ols}} = \boxed{X(X^{\mathsf{T}}X)^{\dagger} X^{\mathsf{T}}} \; \vec{y}$$

projection matrix $\mathbf{\Pi}$

Say $\hat{y}$ is the ols solution, ie,

$$\hat{y} := X \vec{w}_{\text{ols}} = \sum_{i=1}^{d} w_{\text{ols},i} \ddot{x}_i$$

*Thus, $\hat{y}$ is the **orthogonal projection** of y onto the* $\mathrm{span}(\ddot{x}_1, \dots, \ddot{x}_d)$ *!*

$w_{\text{ols}}$ *forms the* **coefficients** *of $\hat{y}$*



residual

$\ddot{x}_2$

$\mathbf{\Pi}(y)$ $\qquad \mathrm{span}(\ddot{x}_1, \dots, \ddot{x}_d)$

$\ddot{x}_1$

# Linear Regression: Statistical Modeling View

Let's assume that data is **generated** from the following process:

- A example $x_i$ is draw independently from the data space **X**

$$x_i \sim \mathcal{D}_X$$

- $y_i^{\text{clean}}$ is computed as ($w . x_i$), from a fixed unknown $w$

$$y_i^{\text{clean}} := w \cdot x_i$$

- $y_i^{\text{clean}}$ is corrupted from by adding independent Gaussian noise $N(0,\sigma^2)$

$$y_i := y_i^{\text{clean}} + \epsilon_i = w \cdot x_i + \epsilon_i \qquad \epsilon_i \sim N(0, \sigma^2)$$

- ($x_i$, $y_i$) is revealed as the $i^{th}$ sample

$$(x_1, y_1), \ldots, (x_n, y_n) =: S$$

# Linear Regression: Statistical Modeling View

How can we determine *w*, from Gaussian noise corrupted observations?

$$S = (x_1, y_1), \ldots, (x_n, y_n)$$

Observation:

$$y_i \sim w \cdot x_i + N(0, \sigma^2) = N(w \cdot x_i, \sigma^2)$$

*How to estimate parameters of a Gaussian?*

*Let's try Maximum Likelihood Estimation!*

parameter

$$\log \mathcal{L}(w|S) \quad = \sum_{i=1}^{n} \log p(y_i|w)$$

$$\propto \sum_{i=1}^{n} \frac{-(w \cdot x_i - y_i)^2}{2\sigma^2}$$
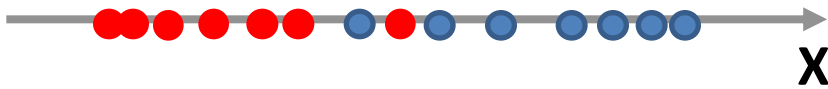
*ignoring terms independent of w*

*optimizing for w yields the same ols result!*

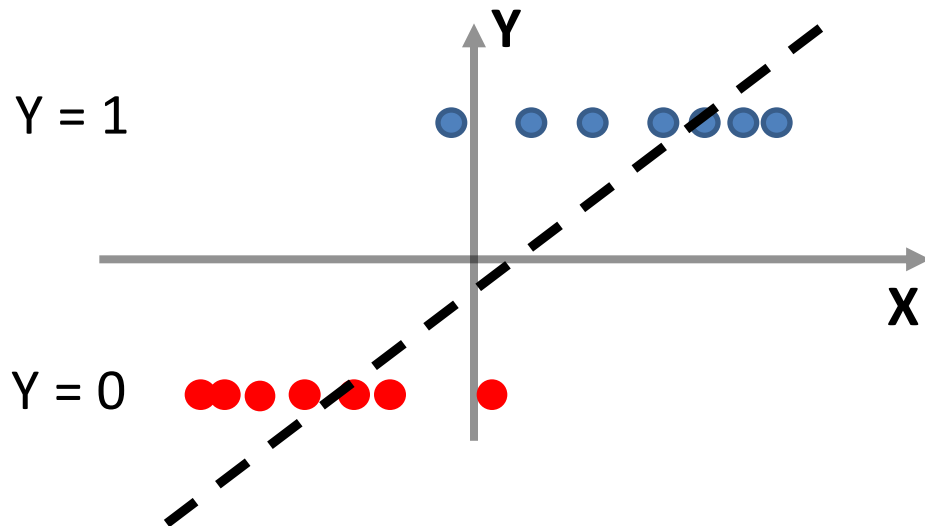*What happens if we model each $y_i$ with indep. noise of different variance?*

# Linear Regression for Classification?

Linear regression seems general, can we use it to derive a binary classifier?
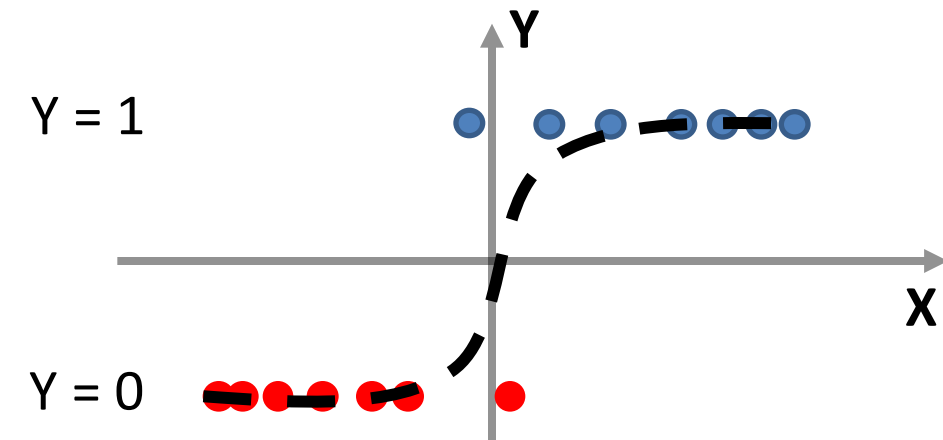
Let's study 1-d data:



*Problem #1: Where is y? for regression.*

*Problem #2: Not really linear!*

*Perhaps it is linear in some **transformed** coordinates?*

# Linear Regression for Classification



*Sigmoid a better model!*

$$\hat{y} = f(x) := \frac{1}{1 + e^{-w \cdot x}}$$

*Binary predictor:* $\operatorname{sign}(2f(x) - 1)$

# Linear Regression for Classification

**Probabilistic Interpretation (of using the sigmoid model):**

For a binary classification problem, given input x, **how likely** is it that it has label 1?

Let this be denoted by P, ie, P is the chance that a given x the associated label y = 1

P = P(Y=1|X=x) ranges between 0 and 1, hence cannot be modelled appropriately via linear regression

How about some other indicator of 'success' (y=1)?

If we look at the 'odds' of getting y=1 (for a given x)

$$\text{odds}(P) := \frac{P}{1-P}$$

*For an event with P=0.9, odds = 9*
*But, for an event P=0.1, odds = 0.11*
*(very asymmetric)*

Consider the "log" of the odds

$$\log(\text{odds}(P)) := \text{logit}(P) := \log\left(\frac{P}{1-P}\right)$$

$$\text{logit}(P) = -\text{logit}(1-P)$$  *Symmetric!*    *Can model logit as a linear function!!*

# Logistic Regression

Model the log-odds or logit with linear function!

Given an input x

modeling assumption

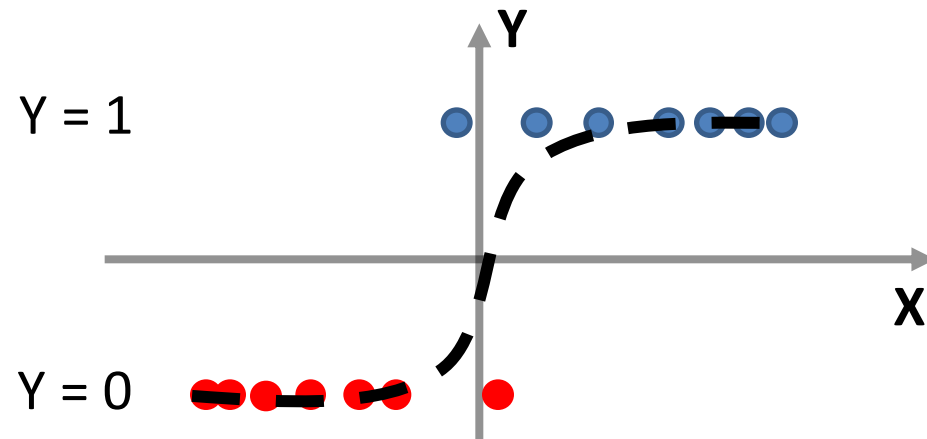$$\text{logit}(P(Y = 1 | X = x)) = \text{logit}(P) = \log\left(\frac{P}{1 - P}\right) = w \cdot x$$

$$\frac{P}{1 - P} = e^{w \cdot x}$$

$$P(Y = 1 | X = x) = \quad P = \frac{e^{w \cdot x}}{1 + e^{w \cdot x}} = \frac{1}{1 + e^{-w \cdot x}} \quad \textit{Sigmoid!}$$

*OK, we have a model, how do we learn the parameters?*

Y = 1

Y = 0

Y

X

# Logistic Regression: Learning Parameters

Given samples $\quad S = (x_1, y_1), \ldots, (x_n, y_n) \qquad (y_i \in \{0,1\}$ binary$)$

$$\mathcal{L}(w|S) = \prod_{i=1}^{n} P((x_i, y_i)|w) \propto \prod_{i=1}^{n} P(y_i|x_i, w)$$

$$= \prod_{i=1}^{n} P(y_i = 1|x_i, w)^{y_i} (1 - P(y_i = 1|x_i, w))^{1-y_i}$$

*(Binomial MLE)*

$$\log \mathcal{L}(w|S) \propto \sum_{i=1}^{n} y_i \log P_{x_i} + (1 - y_i) \log(1 - P_{x_i})$$

$$= \sum_{i=1}^{n} y_i \log \frac{P_{x_i}}{1 - P_{x_i}} + \sum_{i=1}^{n} \log(1 - P_{x_i}) \qquad \textit{Now, use logistic model!}$$

$$= \sum_{i=1}^{n} y_i (w \cdot x_i) + \sum_{i=1}^{n} -\log(1 + e^{w \cdot x_i})$$

*no closed form solution*
*Can use iterative methods like gradient 'ascent' to find the solution*

# Linear Regression: Other Variations

Back to the ordinary least squares (ols):

$$\text{minimize} \ \left\| X\vec{w} - \vec{y} \right\|_2^2$$

$$\vec{w}_{\text{ols}} = (X^\mathsf{T} X)^\dagger X^\mathsf{T} \vec{y}$$

*poorly behaved (due to overfitting) when we have limited data*

how can we incorporate application dependent prior knowledge? e.g.

- perhaps only a few input features dictate/control the outcome  *Lasso regression*

- perhaps multiple features are highly correlated, need to find a stable relationship between the input and the response variable.  *Ridge regression*

# Ridge Regression

Objective

$$\text{minimize} \ \left\| X\vec{w} - \vec{y} \right\|^2 + \lambda \left\| \vec{w} \right\|^2$$

reconstruction error          'regularization' parameter

$$\vec{w}_{\text{ridge}} = (X^\mathsf{T} X + \lambda I)^{-1} X^\mathsf{T} \vec{y}$$
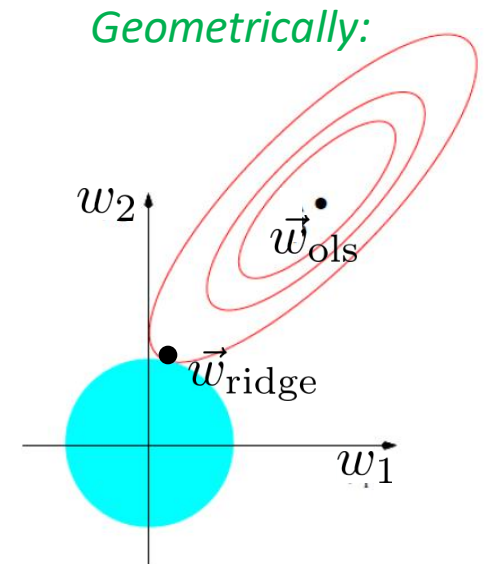
The 'regularization' helps avoid overfitting, and always resulting in a unique solution.

Equivalent to the following optimization problem:

$$\text{minimize} \ \left\| X\vec{w} - \vec{y} \right\|^2$$

$$\text{such that} \ \left\| \vec{w} \right\|^2 \leq B$$

*Why?*

*Geometrically:*

# Lasso Regression

Objective

$$\text{minimize} \ \left\| X\vec{w} - \vec{y} \right\|^2 + \lambda \|\vec{w}\|_1$$

'lasso' penalty

$$\vec{w}_{\text{lasso}} = ?$$   *no closed form solution*

Lasso regularization encourages sparse solutions.
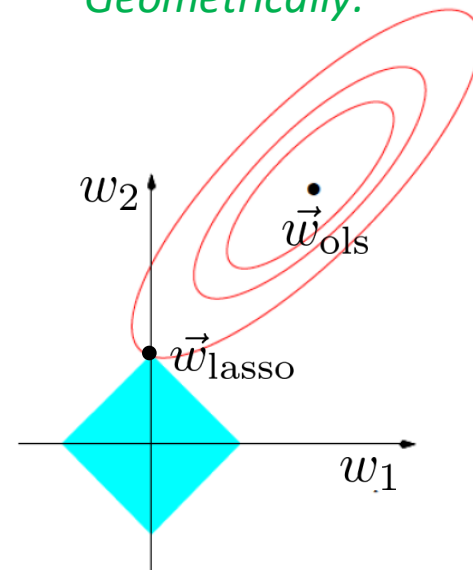
Equivalent to the following optimization problem:

$$\text{minimize} \ \left\| X\vec{w} - \vec{y} \right\|^2$$
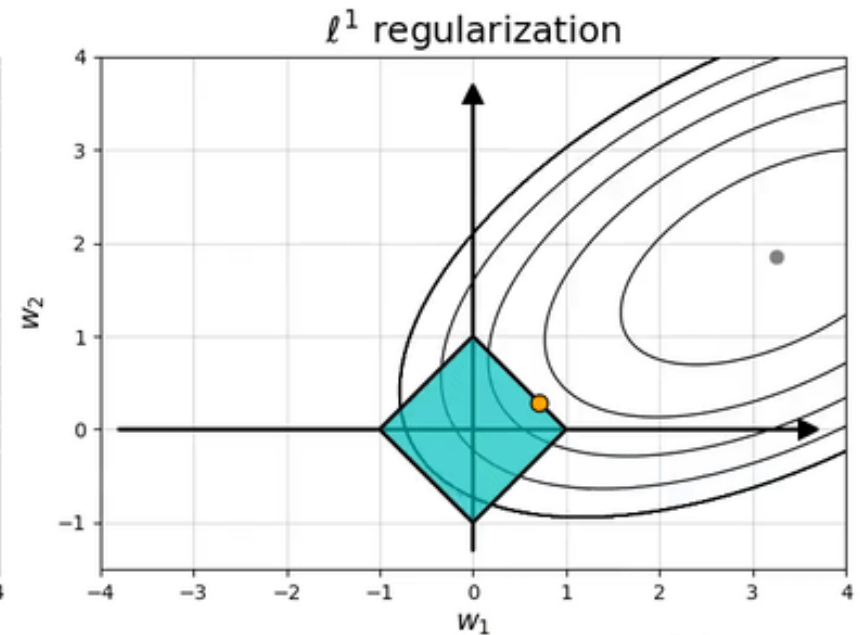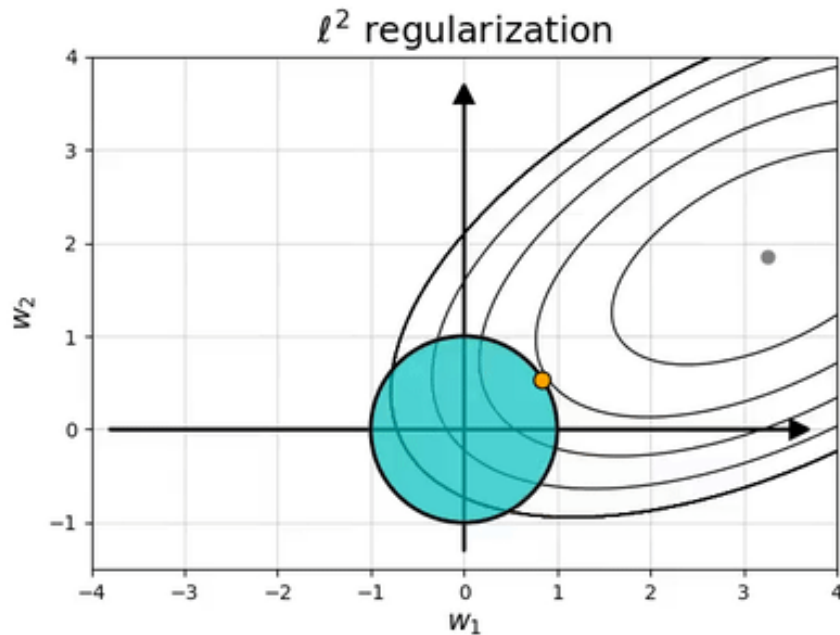$$\text{such that} \ \left\| \vec{w} \right\|_1 \leq B$$

*Why?*

*How can we find the solution?*

*Geometrically:*

# Lasso Regression results in Sparse Solutions



$\ell^1$ induces sparse solutions for least squares

by @itayevron

# What About Optimality?

Linear regression (and variants) is great, but what can we say about the best possible estimate?

*Can we construct an estimator for real outputs that* **parallels** *Bayes classifier for discrete outputs?*

# Optimal L$_2$ Regressor

Best possible regression estimate at *x*:     $f^*(x) := \mathbb{E}\big[Y|X = x\big]$

**Theorem:** for any regression estimate *g(x)*

$$\mathbb{E}_{(x,y)}\big|f^*(x) - y\big|^2 \leq \mathbb{E}_{(x,y)}\big|g(x) - y\big|^2$$

*Similar to Bayes classifier, but for regression.*

*Proof is straightforward…*

# Proof

Consider L$_2$ error of $g(x)$

$$\boxed{f^*(x) := \mathbb{E}\big[Y|X=x\big]}$$

$$\mathbb{E}\big|g(x)-y\big|^2 \;=\; \mathbb{E}\big|g(x)-f^*(x)+f^*(x)-y\big|^2$$

$$= \mathbb{E}\big|g(x)-f^*(x)\big|^2 + \mathbb{E}\big|f^*(x)-y\big|^2 \qquad \textit{Why?}$$

**Cross term:** $\quad 2\mathbb{E}\big[(g(x)-f^*(x))(f^*(x)-y)\big]$

$$= 2\mathbb{E}_x\big[\mathbb{E}_{y|x}\big[(g(x)-f^*(x))(f^*(x)-y)\mid X=x\big]\big]$$

$$= 2\mathbb{E}_x\big[(g(x)-f^*(x))\cdot \mathbb{E}_{y|x}\big[(f^*(x)-y)\mid X=x\big]\big]$$

$$= 2\mathbb{E}_x\big[(g(x)-f^*(x))(f^*(x)-f^*(x))\big] \;=\; 0$$

Therefore $\quad \mathbb{E}\big|g(x)-y\big|^2 = \int_x \big|g(x)-f^*(x)\big|^2 \,\mu(dx) + \mathbb{E}\big|f^*(x)-y\big|^2$

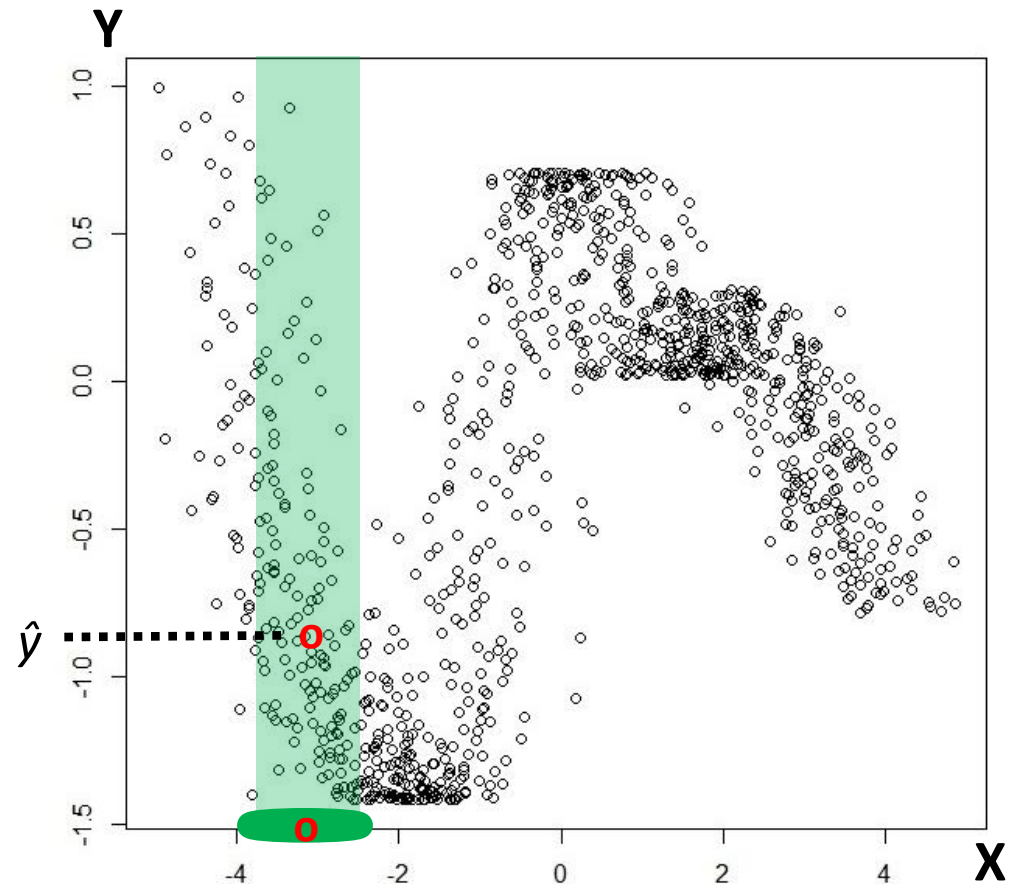*Which is minimized when g(x) = f*(x)!*

■

# Non-parametric Regression

Linear regression (and variants) is great, but what if we don't know parametric form of the relationship between the independent and dependent variables?

How can we predict value of a new test point $x$ *without* model assumptions?

Idea:

$\hat{y} = f(x) =$ *Average estimate **Y** of observed data in a local neighborhood **X** of x!*

# Kernel Regression

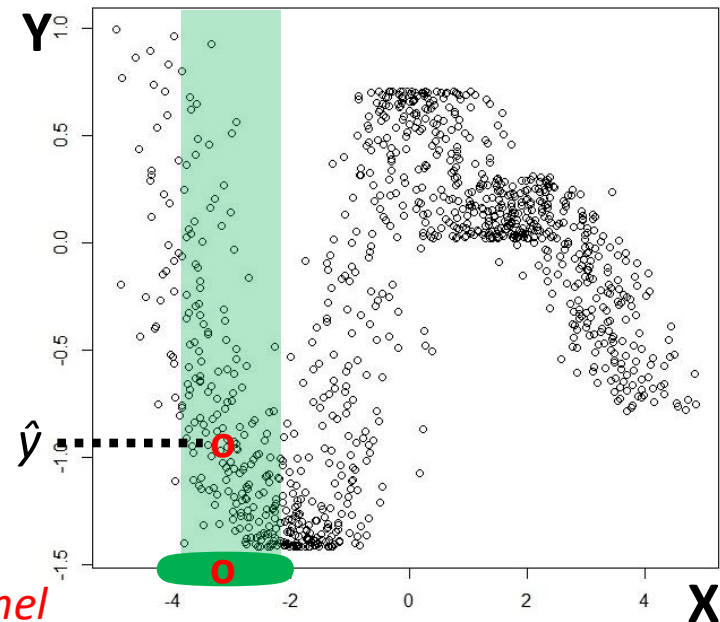$$\hat{y} = \hat{f}_n(x) := \sum_{i=1}^{n} \boxed{w_i(x)} y_i$$

Want weights that emphasize **local** observations

Consider example localization functions:

$$K_h(x, x') = e^{-\|x-x'\|^2/h} \qquad \textit{Gaussian kernel}$$

$$= \mathbf{1}\big[\|x - x'\| \leq h\big] \qquad \textit{Box kernel}$$

$$= \big[1 - (1/h)\|x - x'\|\big]_{+} \qquad \textit{Triangle kernel}$$

Then define:

$$w_i(x) := \frac{K_h(x, x_i)}{\sum_{j=1}^{n} K_h(x, x_j)} \qquad \textit{Weighted average}$$

# Consistency Theorem

Recall: best possible regression estimate at x:   $f^*(x) := \mathbb{E}\big[Y|X=x\big]$

**Theorem:** As $n \to \infty$, $h \to 0$, $hn \to \infty$, then

$$\mathbb{E}_{(\vec{x},y)}\big|\hat{f}_{n,h}(x) - f^*(x)\big|^2 \to 0$$

where   $\hat{f}_{n,h}(x) := \sum_{i=1}^{n} \frac{K_h(x,x_i)}{\sum_{j=1}^{n} K_h(x,x_j)}\, y_i$  is the kernel regressor with most localization kernels.

*Proof is a bit tedious…*

# Proof Sketch

Prove for a fixed x and then integrate over (just like before)

$$\mathbb{E}\left|\hat{f}_{n,h}(x) - f^*(x)\right|^2 = \left[\mathbb{E}\hat{f}_{n,h}(x) - f^*(x)\right]^2 + \mathbb{E}\left[\hat{f}_{n,h}(x) - \mathbb{E}\hat{f}_{n,h}(x)\right]^2$$

squared bias of $\hat{f}_{n,h}$          variance of $\hat{f}_{n,h}$          *Bias-variance decomposition*

Sq. bias   $\approx c_1 h^2$

Variance   $\approx c_2 \dfrac{1}{nh^d}$

Pick    $h \approx n^{-1/2+d}$          $\mathbb{E}\left|\hat{f}_{n,h}(x) - f^*(x)\right|^2 \approx n^{-2/2+d} \to 0$

# Kernel Regression

$$\hat{y} = \hat{f}_n(x) := \sum_{i=1}^{n} \frac{K_h(x, x_i)}{\sum_{j=1}^{n} K_h(x, x_j)} \, y_i$$

Advantages:

- Does not assume any parametric form of the regression function.
- Kernel regression is consistent

Disadvantages:

- Evaluation time complexity:    O($dn$)
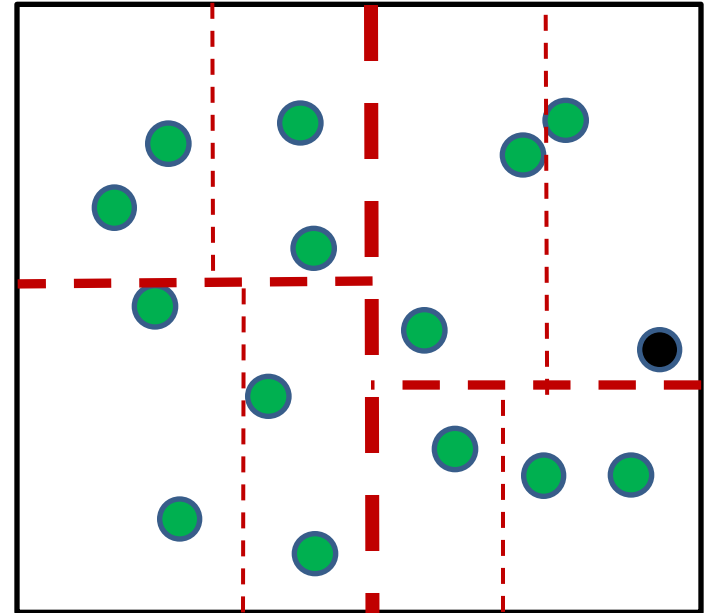- Need to keep all the data around!

*How can we address the shortcomings of kernel regression?*

k-d trees to the rescue!

Idea: partition the data in cells organized in a tree based hierarchy. (just like before)

To return an estimated value, return the average y value in a cell!

# What We Learned…

- Linear Regression

- Parametric vs Nonparametric regression

- Logistic Regression for classification

- Ridge and Lasso Regression

- Kernel Regression

- Consistency of Kernel Regression

- Speeding non-parametric regression with trees

# Questions?

# Next time…

Statistical Theory of Learning!