

Foundations of Machine Learning

Regression

Nov 2021

Vineeth N Balasubramanian



ML Problems

Supervised Learning

Unsupervised Learning

Discrete

classification or
categorization

clustering

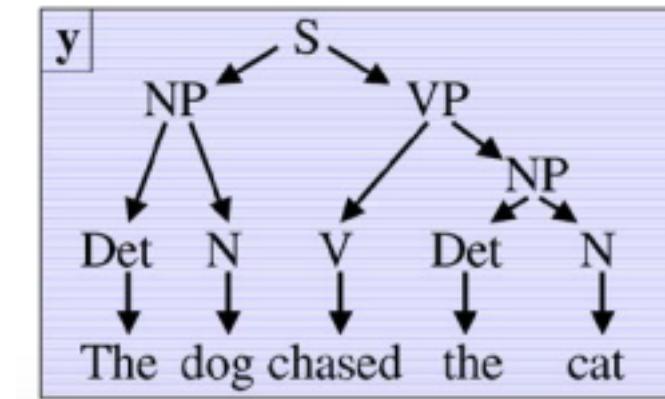
Continuous

regression

dimensionality
reduction

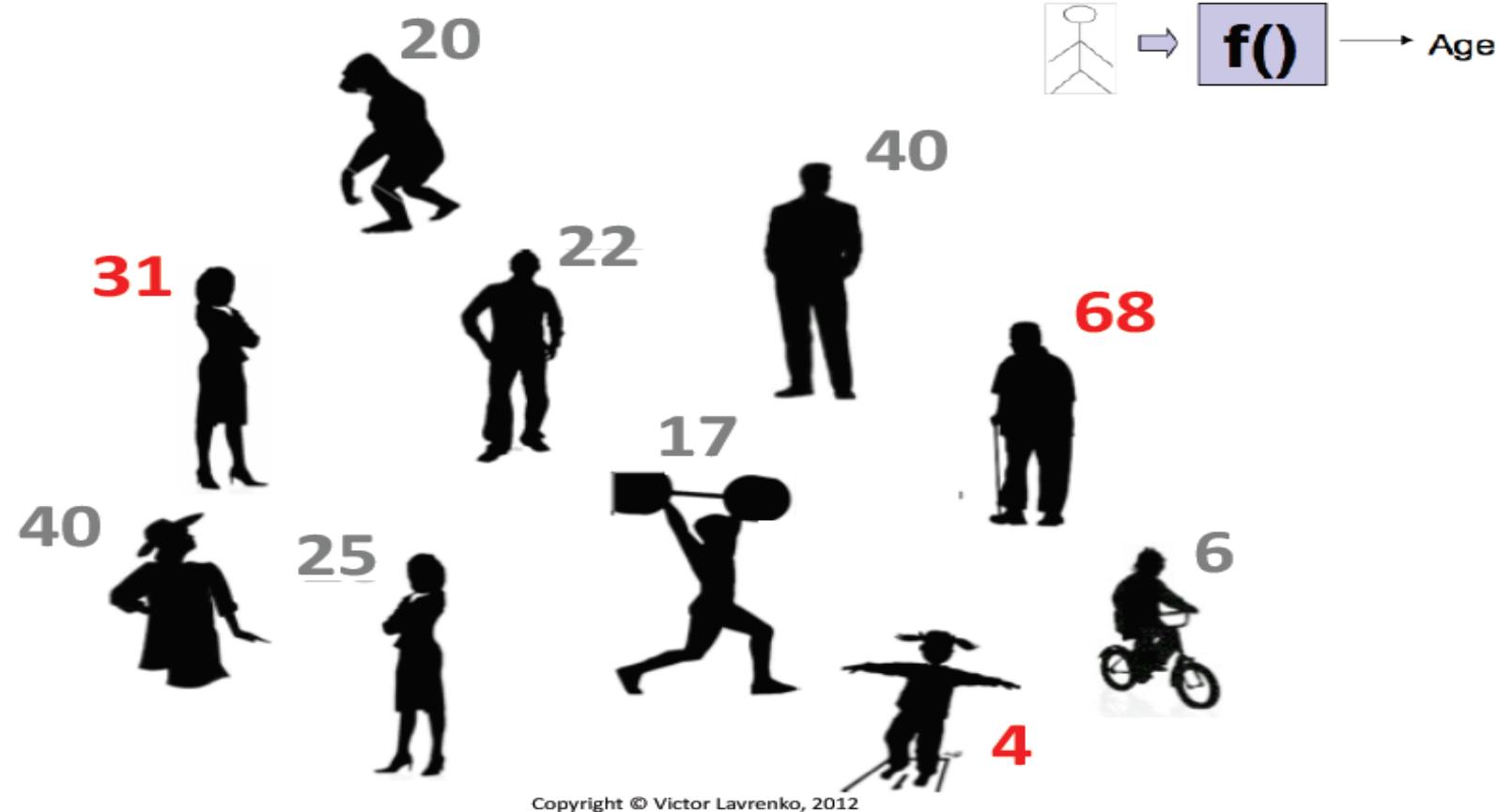
Going beyond Discrete Class Labels

- So far we have focused on classification $f : X \rightarrow \{1, \dots, k\}$
What about other outputs?
- PM2.5 (pollutant) particulate matter exposure estimate:
 - Input: # cars, temperature, etc.
 - Output: 50 ppb
- Pose estimation:
- Sentence structure estimate:



Slide Credit: Nakul Verma, Columbia Univ

Regression (Supervised Learning)

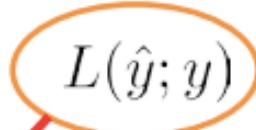


Regression Methods

- Linear Least-Squares Regression
 - Partial Least-Squares
 - Total Least-Squares
 - Ridge Regression, LASSO
- Kernel Regression
- k-NN Regression
- Regression Trees
- Support Vector Regression
- Logistic Regression

Regression Formulation

Given x , want to predict an estimate \hat{y} of y , which minimizes the discrepancy (L) between \hat{y} and y .

Loss 

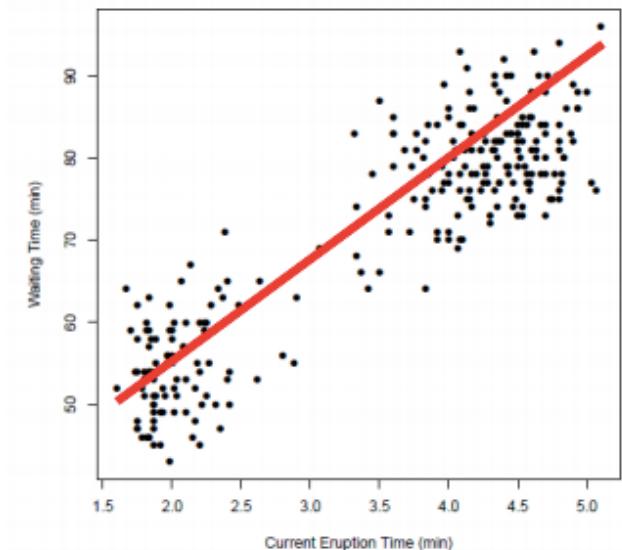
$$\begin{aligned} L(\hat{y}; y) &:= |\hat{y} - y| && \textcolor{red}{\text{Absolute error}} \\ &:= (\hat{y} - y)^2 && \textcolor{red}{\text{Squared error}} \end{aligned}$$

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.

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



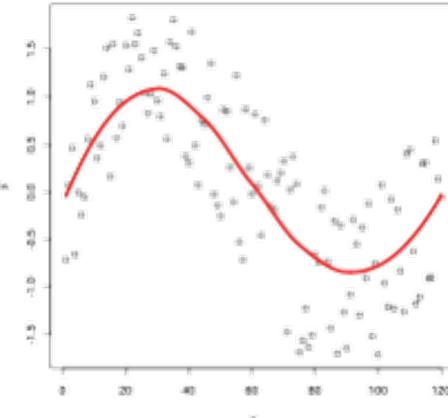
Slide Credit: Nakul Verma, Columbia Univ

Parametric (vs) Non- parametric Regression

If we assume a particular form of the regressor:

Parametric regression

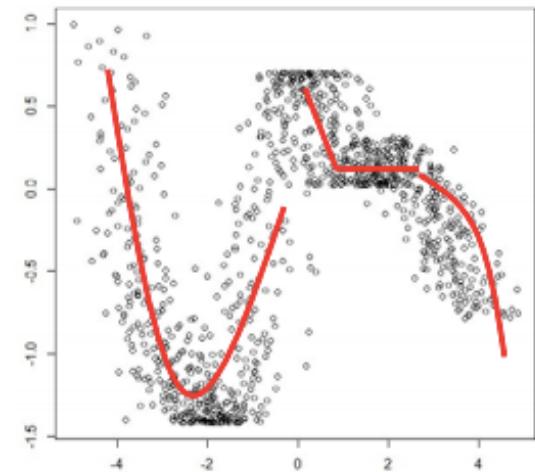
Goal: to learn the parameters which yield the minimum error/loss



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



Slide Credit: Nakul Verma, Columbia Univ

Linear Regression

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

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

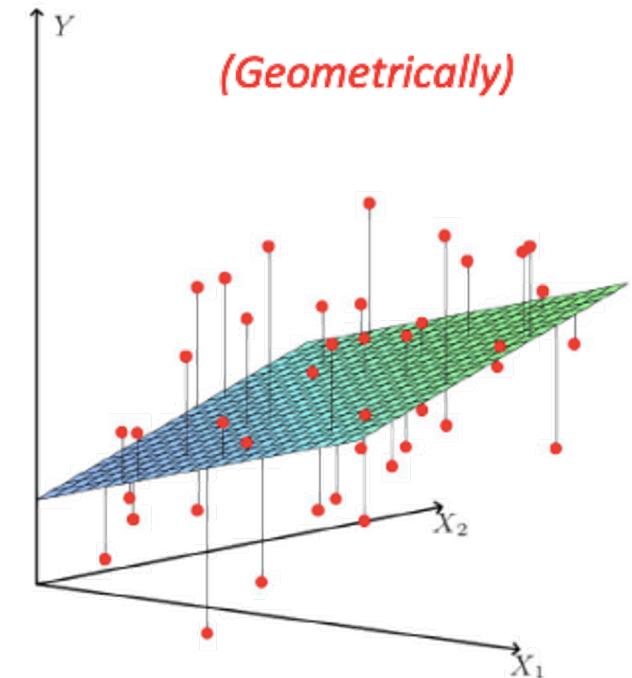
which minimizes the prediction loss over the population.

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

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

$$\begin{aligned} & \arg \min_w \frac{1}{n} \sum_{i=1}^n [L(\vec{w} \cdot \vec{x}_i, y_i)] \\ &= \arg \min_w \frac{1}{n} \sum_{i=1}^n (\vec{w} \cdot \vec{x}_i - y_i)^2 \end{aligned}$$

for squared error



Slide Credit: Nakul Verma, Columbia Univ

Linear Regression

Linear predictor with squared loss:

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

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

$$= \arg \min_w \|X\vec{w} - \vec{y}\|_2^2$$

Unconstrained problem!

Can take the gradient and examine the stationary points!

Why need not check the second order conditions?

Slide Credit: Nakul Verma, Columbia Univ



Foundations of Machine Learning

Linear Regression

Best fitting w :

$$\frac{\partial}{\partial \vec{w}} \|X\vec{w} - \vec{y}\|^2 = 2X^T(X\vec{w} - \vec{y})$$

$$X^T X \vec{w} = X^T \vec{y}$$

At a stationary point

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

Pseudo-inverse

Also called the Ordinary Least Squares (OLS)

The solution is unique and stable when $X^T X$ is invertible

What is the interpretation of this solution?

Slide Credit: Nakul Verma, Columbia Univ



Linear Regression: Geometric Interpretation

Consider the **column space** view of data \mathbf{X} :

$$\left(\begin{array}{c|c} \cdots & \mathbf{x}_1 \cdots \\ \cdots & \mathbf{x}_i \cdots \\ \cdots & \mathbf{x}_n \cdots \end{array} \right) \quad \ddot{x}_1, \dots, \ddot{x}_d \in \mathbf{R}^n$$

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

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

$$\hat{y} = \mathbf{X}\vec{w}_{\text{ols}} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^\dagger \mathbf{X}^\top \vec{y}$$

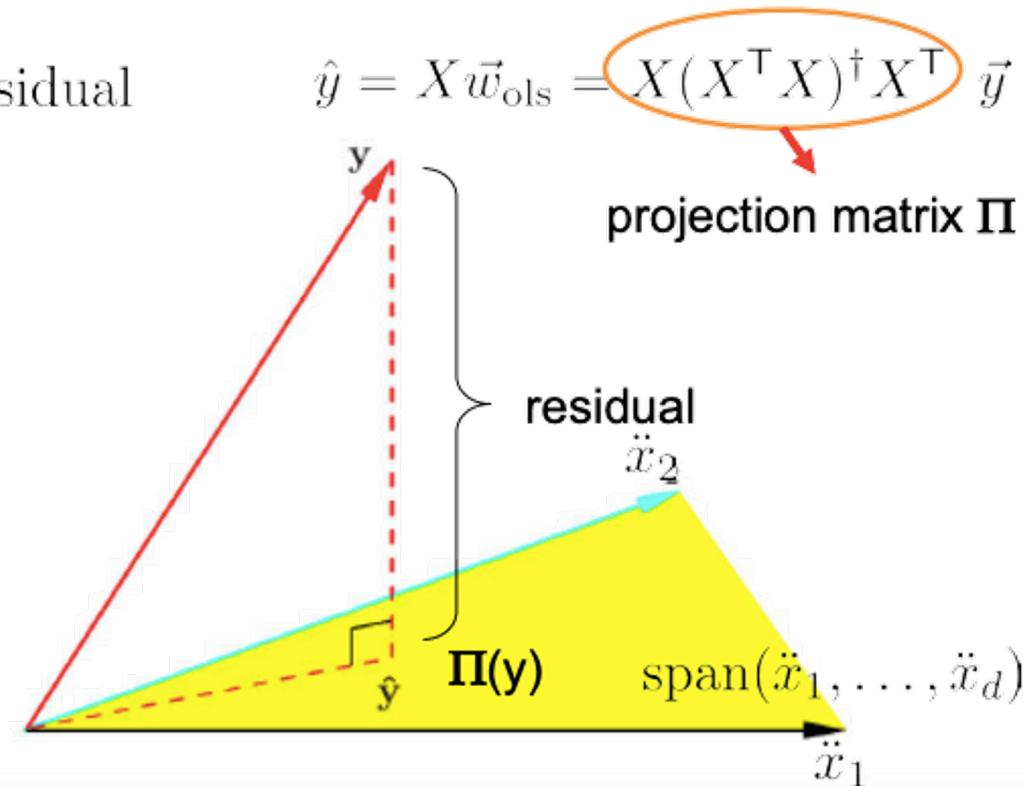
projection matrix Π

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

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

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

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



Slide Credit: Nakul Verma, Columbia Univ

Linear Regresssion: 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 \mathbf{X}

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

- y_{clean} is computed as $(w \cdot x_i)$, from a fixed unknown w

$$y_{\text{clean}} := w \cdot x_i$$

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

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

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

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

Linear Regression: Statistical Modeling View

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

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

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?*

Learning Paradigms

Statistical modeling approach:

Labeled training data
(n examples from data)

$(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots (\vec{x}_n, y_n)$

drawn **independently** from
a fixed underlying distribution
(also called the *i.i.d.* assumption)



select \hat{f} from...?

from a pool of **models** \mathcal{F}
that **maximizes**
label agreement of the
training data

How to select $\hat{f} \in \mathcal{F}$?

- Maximum likelihood (best fits the data)
- Maximum a posteriori (best fits the data but incorporates prior assumptions)
- Optimization of ‘loss’ criterion (best discriminates the labels)
- ...

Slide Credit: Nakul Verma, Columbia Univ

Maximum Likelihood Estimation

Given some data $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathcal{X}$ i.i.d. (Let's forget about the labels for now)

Say we have a model class $\mathcal{P} = \{p_\theta \mid \theta \in \Theta\}$

ie, each model p can be described by a set of parameters θ

find the parameter settings θ that **best fits** the data.

If each model p , is a **probability model** then we can find the best fitting probability model via the **likelihood estimation!**

Likelihood $\mathcal{L}(\theta|X) := P(X|\theta) = P(\vec{x}_1, \dots, \vec{x}_n|\theta) \stackrel{i.i.d.}{=} \prod_{i=1}^n P(\vec{x}_i|\theta) = \prod_{i=1}^n p_\theta(\vec{x}_i)$

Interpretation: How **probable** (or how likely) is the data given the model p_θ ?

Parameter setting θ that maximizes $\mathcal{L}(\theta|X)$

$$\arg \max_{\theta} \mathcal{L}(\theta|X) = \arg \max_{\theta} \prod_{i=1}^n p_\theta(\vec{x}_i)$$

Slide Credit: Nakul Verma, Columbia Univ

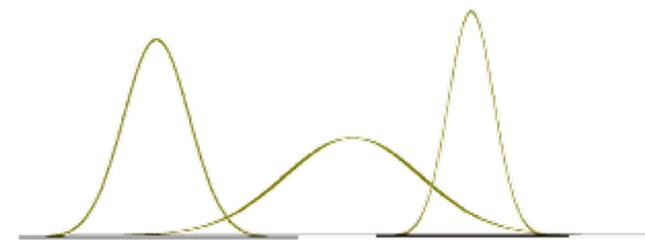
MLE Example

Fitting a statistical probability model to heights of females

Height data (in inches): $60, 62, 53, 58, \dots \in \mathbb{R}$

$$x_1, x_2, \dots, x_n \in \mathcal{X}$$

Model class: Gaussian models in **R**



$$p_{\theta}(x) = p_{\{\mu, \sigma^2\}}(x) := \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

μ = mean parameter
 σ^2 = variance parameter > 0

So, what is the MLE for the given data X ?

Slide Credit: Nakul Verma, Columbia Univ

MLE Example

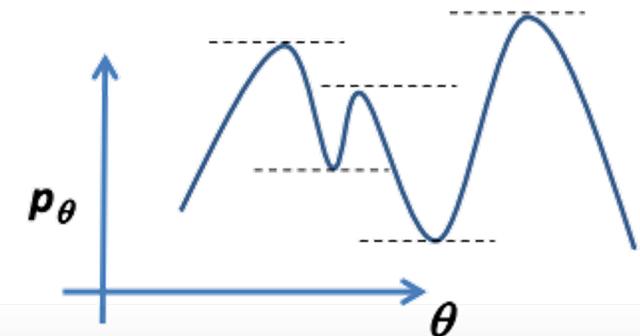
Height data (in inches): $x_1, x_2, \dots, x_n \in \mathcal{X} = \mathbf{R}$

Model class: Gaussian models in \mathbf{R} $p_{\{\mu, \sigma^2\}}(x) := \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

MLE: $\arg \max_{\theta} \mathcal{L}(\theta|X) = \arg \max_{\mu, \sigma^2} \prod_{i=1}^n p_{\{\mu, \sigma^2\}}(x_i)$ **Good luck!**

Trick #1: $\arg \max_{\theta} \mathcal{L}(\theta|X) = \arg \max_{\theta} \log \mathcal{L}(\theta|X)$ **"Log" likelihood**

Trick #2: finding max (or other extreme values) of a function is simply analyzing the '**stationary points**' of a function. That is, values at which the **derivative** of the function is zero !



Slide Credit: Nakul Verma, Columbia Univ

MLE Example

Let's calculate the best fitting $\theta = \{\mu, \sigma^2\}$

$$\begin{aligned}
 \arg \max_{\theta} \mathcal{L}(\theta | X) &= \arg \max_{\theta} \log \mathcal{L}(\theta | X) && \text{"Log likelihood"} \\
 &= \arg \max_{\mu, \sigma^2} \log \left(\prod_{i=1}^n p_{\{\mu, \sigma^2\}}(x_i) \right) && \text{i.i.d.} \\
 &= \arg \max_{\mu, \sigma^2} \sum_{i=1}^n \log \left(p_{\{\mu, \sigma^2\}}(x_i) \right) \\
 &= \arg \max_{\mu, \sigma^2} \sum_{i=1}^n \left[-\frac{1}{2} \underbrace{\log(2\pi\sigma^2)}_{g_i(\mu, \sigma^2)} - \frac{(x_i - \mu)^2}{2\sigma^2} \right]
 \end{aligned}$$

Maximizing μ : $0 = \nabla_{\mu} \left(\sum_{i=1}^n g_i(\mu, \sigma^2) \right) \implies \mu_{\text{ML}} = \frac{1}{n} \sum_{i=1}^n x_i$

Maximizing σ^2 : $\sigma_{\text{ML}}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$

Slide Credit: Nakul Verma, Columbia Univ



MLE Example

Let's calculate the best fitting $\theta = \{\mu, \sigma^2\}$

$$\begin{aligned}\arg \max_{\theta} \mathcal{L}(\theta | X) &= \arg \max_{\theta} \log \mathcal{L}(\theta | X) && \text{"Log likelihood"} \\ &= \arg \max_{\mu, \sigma^2} \log \left(\prod_{i=1}^n p_{\{\mu, \sigma^2\}}(x_i) \right) && \text{i.i.d.} \\ &= \arg \max_{\mu, \sigma^2} \sum_{i=1}^n \log \left(p_{\{\mu, \sigma^2\}}(x_i) \right) \\ &= \arg \max_{\mu, \sigma^2} \sum_{i=1}^n \left[-\frac{1}{2} \underbrace{\log(2\pi\sigma^2)}_{g_i(\mu, \sigma^2)} - \frac{(x_i - \mu)^2}{2\sigma^2} \right]\end{aligned}$$

Maximizing μ : $0 = \nabla_{\mu} \left(\sum_{i=1}^n g_i(\mu, \sigma^2) \right) \implies \mu_{\text{ML}} = \frac{1}{n} \sum_{i=1}^n x_i$

Maximizing σ^2 : $\sigma_{\text{ML}}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$

Slide Credit: Nakul Verma, Columbia Univ



Back to Linear Regression

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

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

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?*

parameter


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

*Let's try Maximum
Likelihood Estimation!*

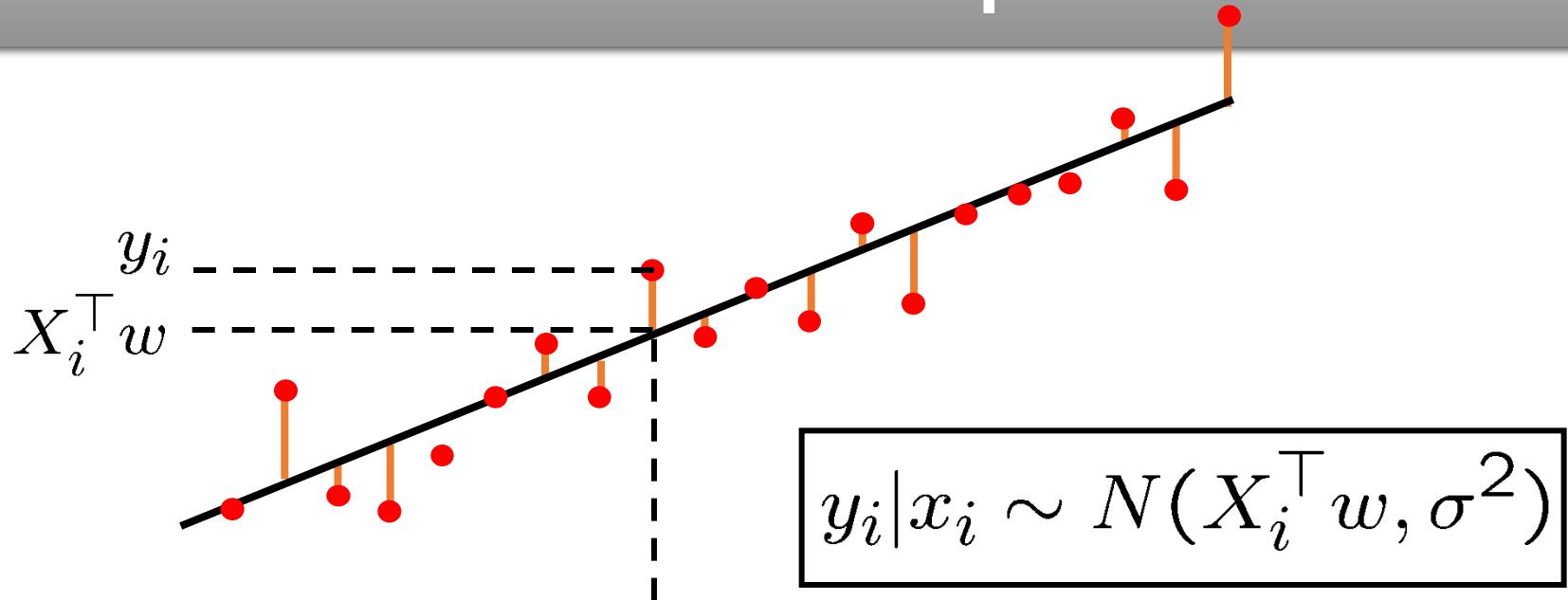
$$\propto \sum_{i=1}^n \frac{-(w \cdot x_i - y_i)^2}{2\sigma^2}$$

*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?**

Probabilistic Interpretation



Likelihood $L = \prod_i \exp -\frac{1}{2\sigma^2} (X_i^\top w - y_i)^2 = \exp -\frac{1}{2\sigma^2} \sum_i (X_i^\top w - y_i)^2$

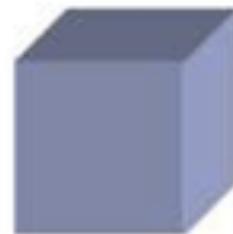
$\underset{w}{\operatorname{argmax}} L = \underset{w}{\operatorname{argmin}} E$

Regularized Least-Squared Regression

- Complex models (lots of parameters) often prone to overfitting.
- Overfitting can be reduced by imposing a constraint on the overall magnitude of the parameters.
- Two common types of regularization in linear regression:
 - **L₂ regularization (a.k.a. ridge regression):** Find \mathbf{w} which minimizes:
$$\sum_{j=1}^N (y_j - \sum_{i=0}^d w_i \cdot x_i)^2 + \lambda \sum_{i=1}^d w_i^2$$
• λ is the regularization parameter: bigger λ imposes more constraint
 - **L₁ regularization (a.k.a. lasso):** Find \mathbf{w} which minimizes:
$$\sum_{j=1}^N (y_j - \sum_{i=0}^d w_i \cdot x_i)^2 + \lambda \sum_{i=1}^d |w_i|$$

Regularized Least-Squared Regression

- Understanding norms and regularization
 - L_p -norms:



$p = \infty$



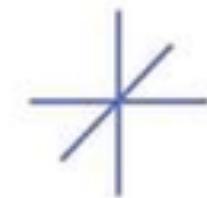
$p = 2$



$p = 1$



$0 < p < 1$



$p = 0$

Solving Ridge Regression

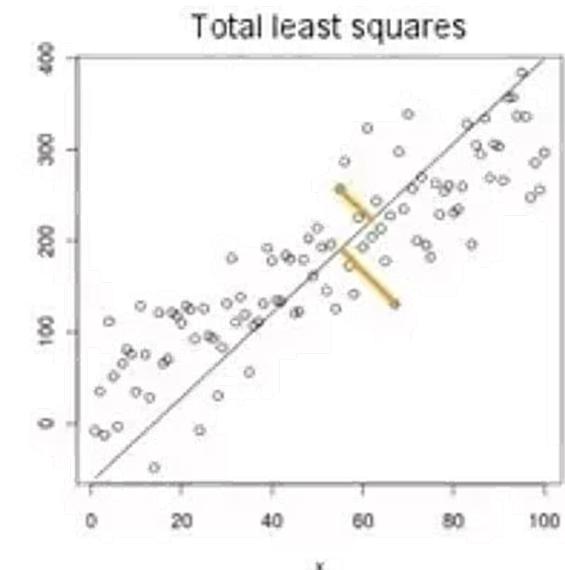
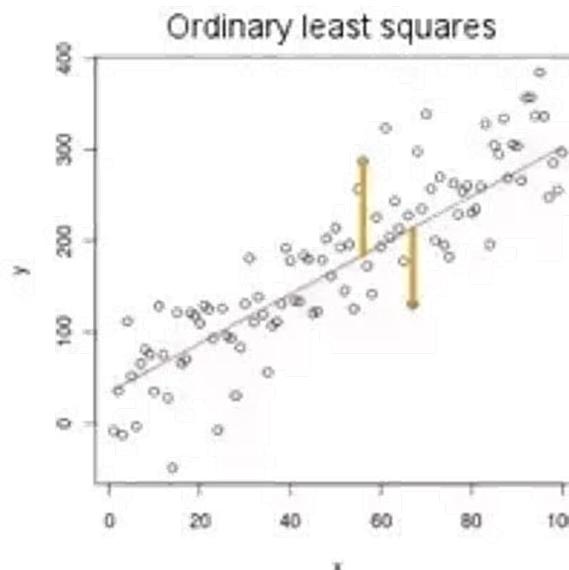
- Minimizing $\lambda\|\mathbf{w}\|^2 + \|\mathbf{y} - \mathbf{Xw}\|^2$ over \mathbf{w} , we get

$$(\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_n)\mathbf{w} = \mathbf{X}'\mathbf{y}$$

- Solution: $\mathbf{w} = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I})^{-1} \mathbf{X}'\mathbf{y}$
- How is this different from the solution for OLS (Ordinary Least-Squares Regression)?
 - Inverse always exists for any $\lambda > 0$.

Total Least Squares and Partial Least Squares

- **Total Least Squares:** Model error in output and input
- **Partial Least Squares:**
 - Seeks to address the correlation between predictor variables in its model
 - Also called “Projection to Latent Structures” (PLS)

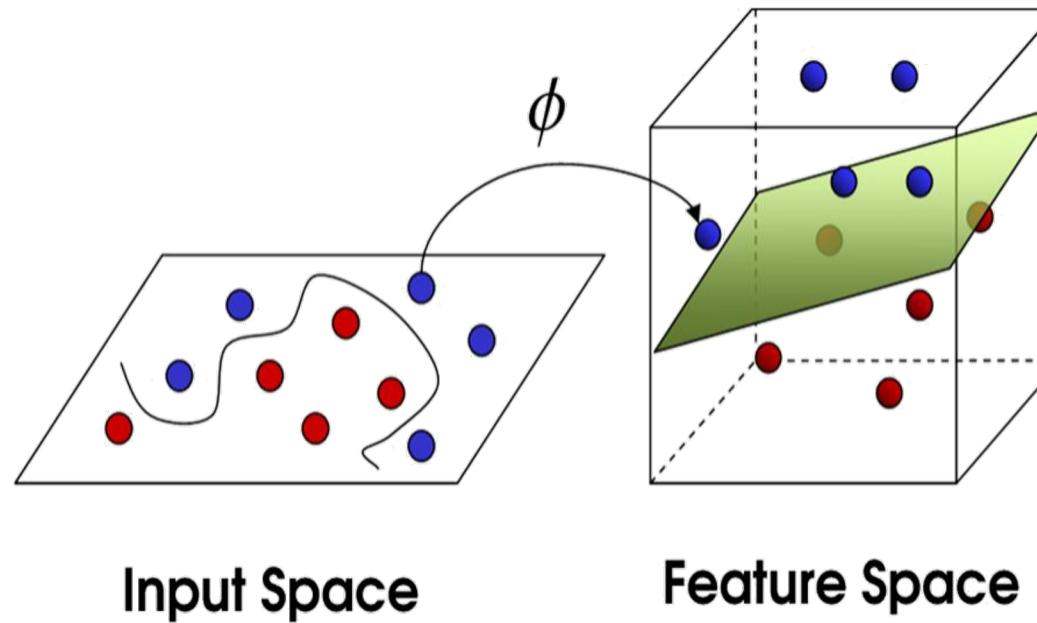


Regression Methods

- Linear Least-Squares Regression
 - Partial Least-Squares
 - Total Least-Squares
 - Ridge Regression, LASSO
- Kernel Regression
- k-NN Regression
- Regression Trees
- Support Vector Regression
- Logistic Regression

Non-Linear Regression

- Recall: “kernel trick”
- **Key Idea:** Map data to higher dimensional space (feature space) and perform linear regression in embedded space



Non-Linear/Kernel Regression

- Alternative view to ridge regression solution:

$$\mathbf{w} = (\mathbf{X}'\mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}'\mathbf{y}$$
Solving is $O(d^3)$

$$(\mathbf{X}'\mathbf{X} + \lambda \mathbf{I})\mathbf{w} = \mathbf{X}'\mathbf{y} \Rightarrow \mathbf{w} = \lambda^{-1} (\mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\mathbf{w})$$

$$\Rightarrow \mathbf{w} = \lambda^{-1} \mathbf{X}' (\mathbf{y} - \mathbf{X}\mathbf{w}) = \mathbf{X}'\mathbf{a}$$

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

$$\Rightarrow \lambda\mathbf{a} = (\mathbf{y} - \mathbf{X}\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{X}'\mathbf{a})$$

$$\Rightarrow \mathbf{X}\mathbf{X}'\mathbf{a} + \lambda\mathbf{a} = \mathbf{y}$$

$$\Rightarrow \mathbf{a} = (\mathbf{G} + \lambda \mathbf{I}_\ell)^{-1} \mathbf{y} \text{ where } \mathbf{G} = \mathbf{X}\mathbf{X}'$$
Solving is $O(n^3)$

Do you spot
anything
interesting?

Non-Linear/Kernel Regression

- To predict new point:

$$g(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle = \left\langle \sum_{i=1}^d \alpha_i \mathbf{x}_i, \mathbf{x} \right\rangle = \mathbf{y}' (\mathbf{G} + \lambda \mathbf{I})^{-1} \mathbf{z}$$

where $\mathbf{z} = \langle \mathbf{x}_i, \mathbf{x} \rangle$

- Need to only compute \mathbf{G} , the Gram Matrix (or the inner products between data points)

$$\mathbf{G} = \mathbf{X}\mathbf{X}' \quad G_{ij} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

Kernel Ridge Regression

- To predict new point:

$$g(\phi(\mathbf{x})) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle = \left\langle \sum_{i=1}^{\ell} \alpha_i \phi(\mathbf{x}_i), \phi(\mathbf{x}) \right\rangle = \mathbf{y}' (\mathbf{G} + \lambda \mathbf{I})^{-1} \mathbf{z}$$

where $\mathbf{z} = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle$

- Use kernel to compute inner products

$$\mathbf{G} = \phi(\mathbf{X}) \phi(\mathbf{X})' \quad G_{ij} = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = K(\mathbf{x}_i, \mathbf{x}_j)$$

Regression Methods

- Linear Least-Squares Regression
 - Partial Least-Squares
 - Total Least-Squares
 - Ridge Regression, LASSO
- Kernel Regression
- **k-NN Regression**
- Regression Trees
- Support Vector Regression
- Logistic Regression

k-NN Regression

- Calculate the mean value of the k nearest training examples rather than calculate their most common value

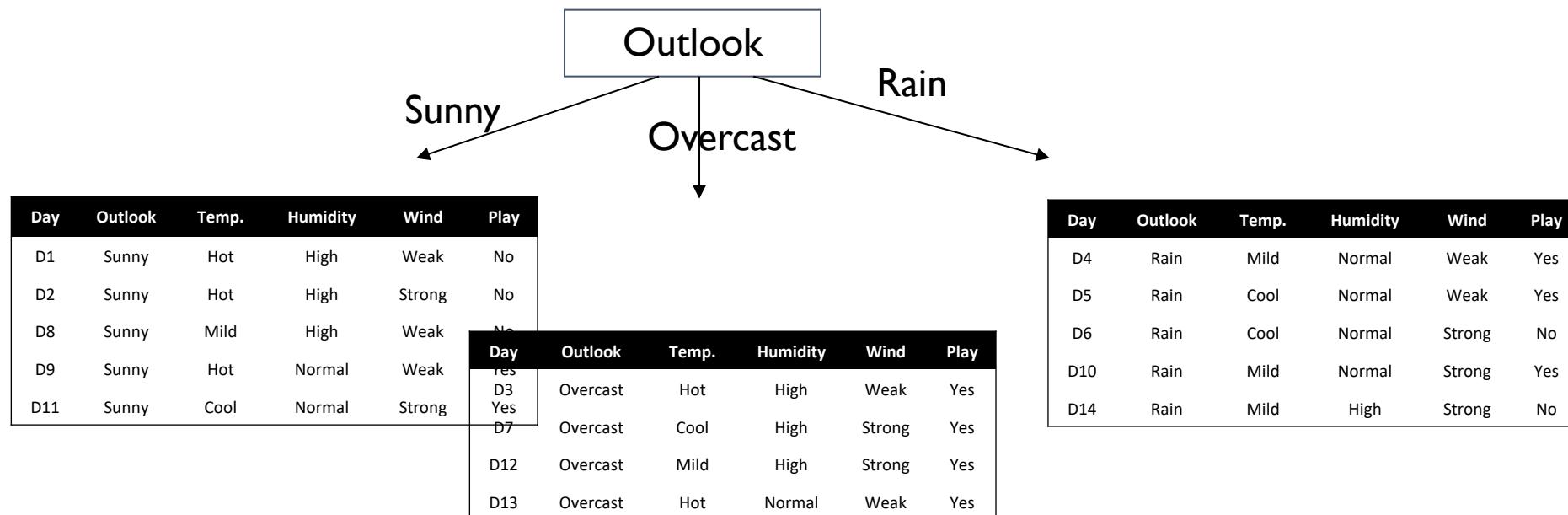
$$f : \Re^d \rightarrow \Re \quad \hat{f}(x_q) \leftarrow \frac{\sum_{i=1}^k f(x_i)}{k}$$

Regression Methods

- Linear Least-Squares Regression
 - Partial Least-Squares
 - Total Least-Squares
 - Ridge Regression, LASSO
- Kernel Regression
- k-NN Regression
- **Regression Trees**
- Support Vector Regression
- Logistic Regression

Recall: Decision Trees

- Choose « best » attribute
- Split the learning sample
- Proceed recursively until each object is correctly classified



Recall: Decision Trees

- The “best” split is the split that maximizes the expected reduction of impurity

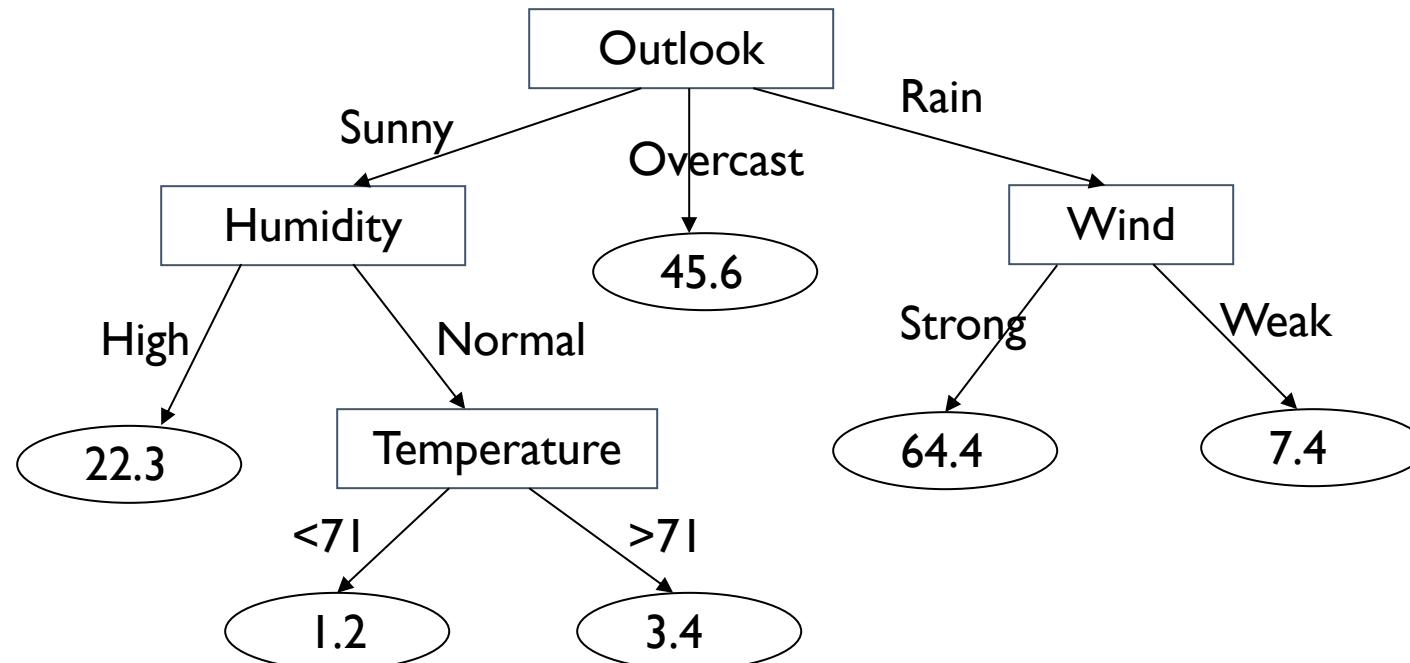
$$\Delta I(LS, A) = I(LS) - \sum_a \frac{|LS_a|}{|LS|} I(LS_a)$$

where LS_a is the subset of records from LS (dataset) such that $A=a$

- Example of impurity measure:
 - *Shannon entropy*: $I(LS) = -\sum_j p_j \log p_j$
 - If two classes, $p_1 = 1 - p_2$
 - The reduction of entropy is called the **information gain**

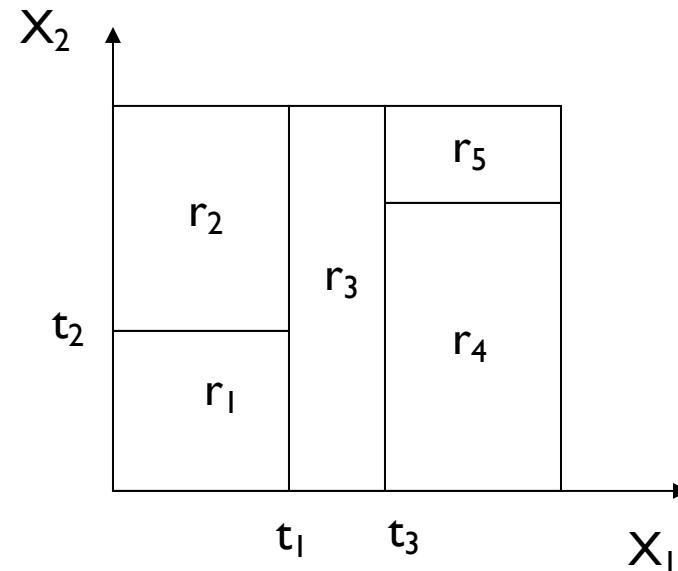
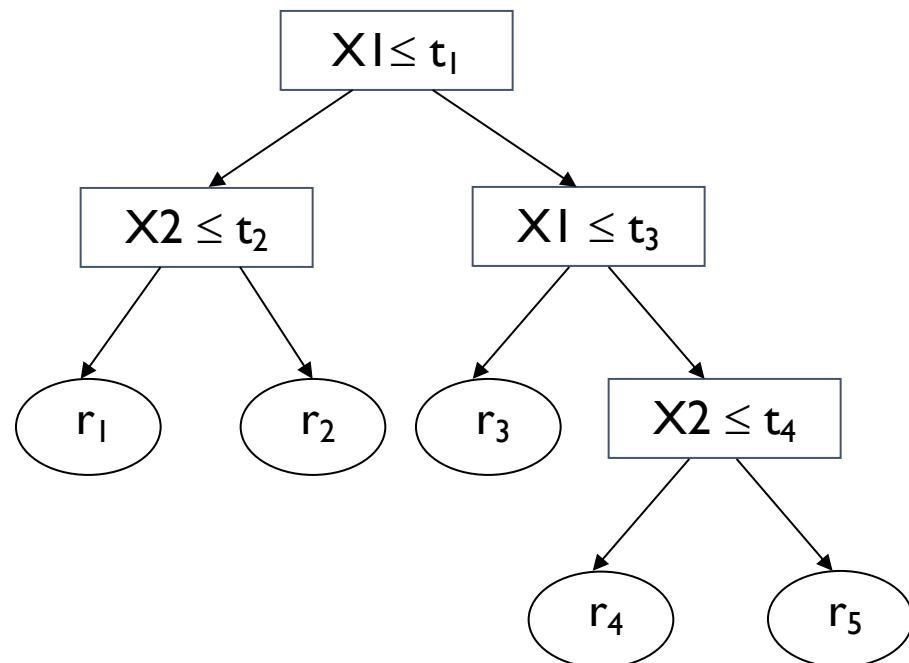
Regression Trees

- Tree for regression: exactly the same model but with a number in each leaf instead of a class



Regression Trees

- A regression tree is a piecewise constant function of the input attributes



Growing Regression Trees

- To minimize the square error on the learning sample, the prediction at a leaf is the average output of the learning cases reaching that leaf
- Impurity of a sample is defined by the variance of the output in that sample:

$$I(LS) = \text{var}_{y|LS}\{y\} = E_{y|LS}\{(y - E_{y|LS}\{y\})^2\}$$

- The best split is the one that reduces the most variance:

$$\Delta I(LS, A) = \text{var}_{y|LS}\{y\} - \sum_a \frac{|LS_a|}{|LS|} \text{var}_{y|LS_a}\{y\}$$

Regression Tree Pruning

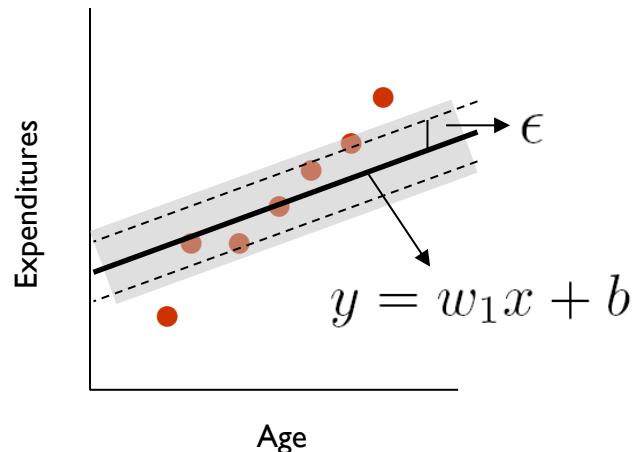
- Exactly the same algorithms apply: pre-pruning and post-pruning.
- In practice, pruning is more important in regression because full trees are much more complex
 - Each data instance can have a different output value and hence the full tree has as many leaves as there are training instances

Regression Methods

- Linear Least-Squares Regression
 - Partial Least-Squares
 - Total Least-Squares
 - Ridge Regression, LASSO
- Kernel Regression
- k-NN Regression
- Regression Trees
- Support Vector Regression
- Logistic Regression

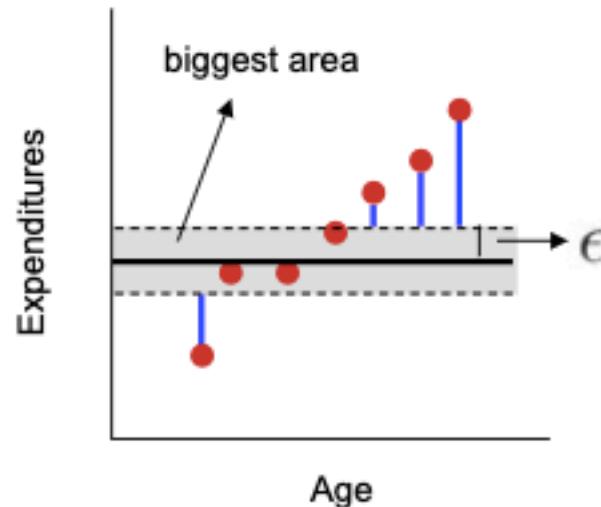
Support Vector Regression

- Given training data $\{x_i, y_i\}_{i=1}^n$
- Find: w_1 and b ,
such that $y = w_1x + b$ optimally describes the data:



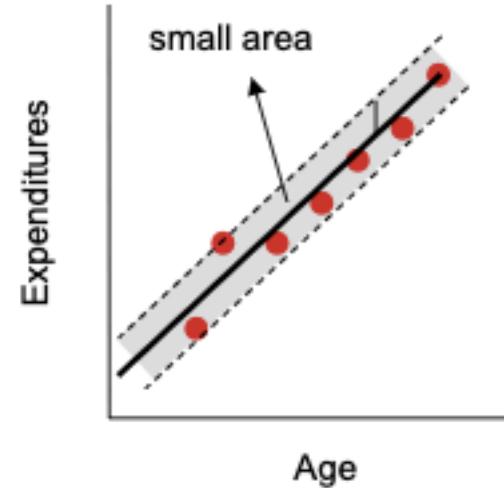
Slide Credit: Georgi Nalbantov, Erasmus Research Institute of Management

Support Vector Regression

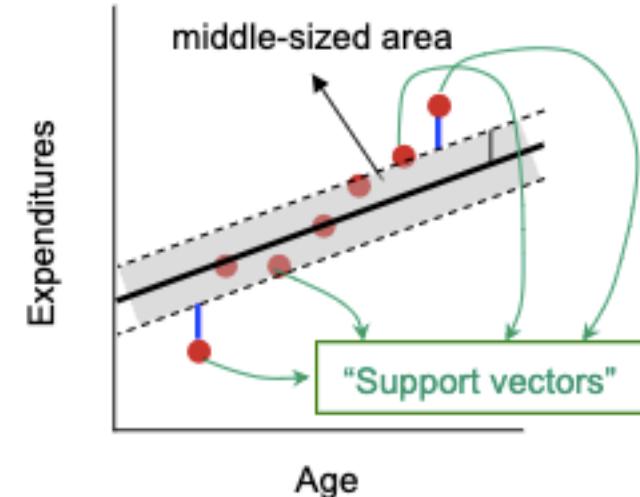


“Lazy case”

(underfitting)



“Suspiciously
smart case”
(overfitting)

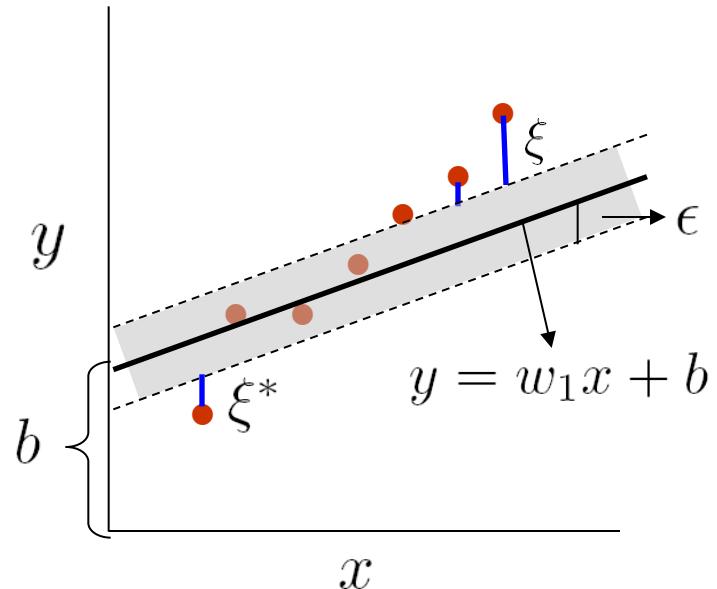


“Compromise case”, SVR
(good generalizability)

- **The thinner the “tube”, the more complex the model**

Slide Credit: Georgi Nalbantov, Erasmus Research Institute of Management

Support Vector Regression



$$|w_1| \text{ vs. } \sum_i (\xi_i + \xi_i^*)$$

Complexity

Sum of errors

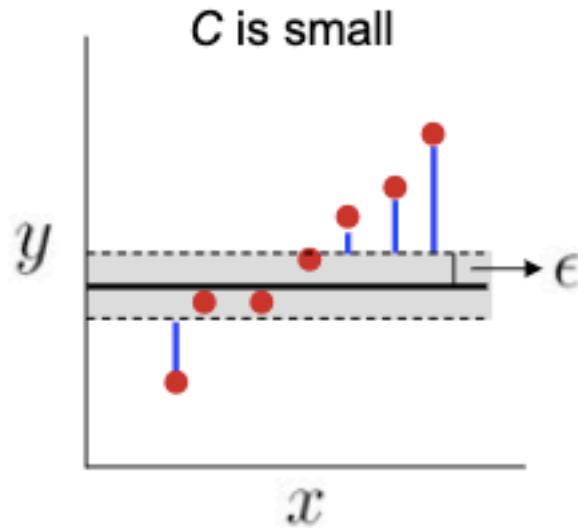
$$\min_{w_1, b, \xi_i, \xi_i^*} \frac{1}{2} w_1^2 + C \sum_i (\xi_i + \xi_i^*)$$

Case I: $w_1 \downarrow \rightarrow$ "tube" $\uparrow \rightarrow$ complexity $\downarrow \rightarrow$ $\sum_i (\xi_i + \xi_i^*) \uparrow$

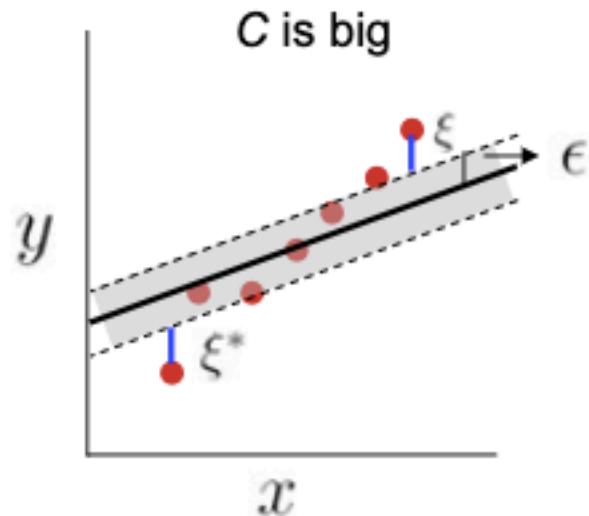
Case II: $w_1 \uparrow \rightarrow$ "tube" $\downarrow \rightarrow$ complexity $\uparrow \rightarrow$ $\sum_i (\xi_i + \xi_i^*) \downarrow$

Slide Credit: Georgi Nalbantov, Erasmus Research Institute of Management

Support Vector Regression



- The role of C

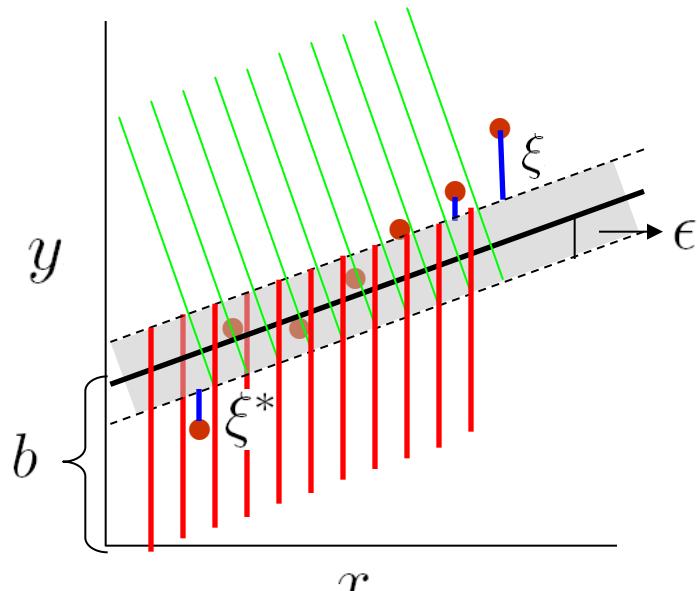


Case I: $w_1 \downarrow \rightarrow$ "tube" $\uparrow \rightarrow$ complexity $\downarrow \rightarrow \sum_i(\xi_i + \xi_i^*) \uparrow$

Case II: $w_1 \uparrow \rightarrow$ "tube" $\downarrow \rightarrow$ complexity $\uparrow \rightarrow \sum_i(\xi_i + \xi_i^*) \downarrow$

Slide Credit: Georgi Nalbantov, Erasmus Research Institute of Management

Support Vector Regression



$$y = w_1x + b$$

$$\min_{w_1, b, \xi_i, \xi_i^*} \frac{1}{2} w_1^2 + C \sum_i (\xi_i + \xi_i^*)$$

Subject to:

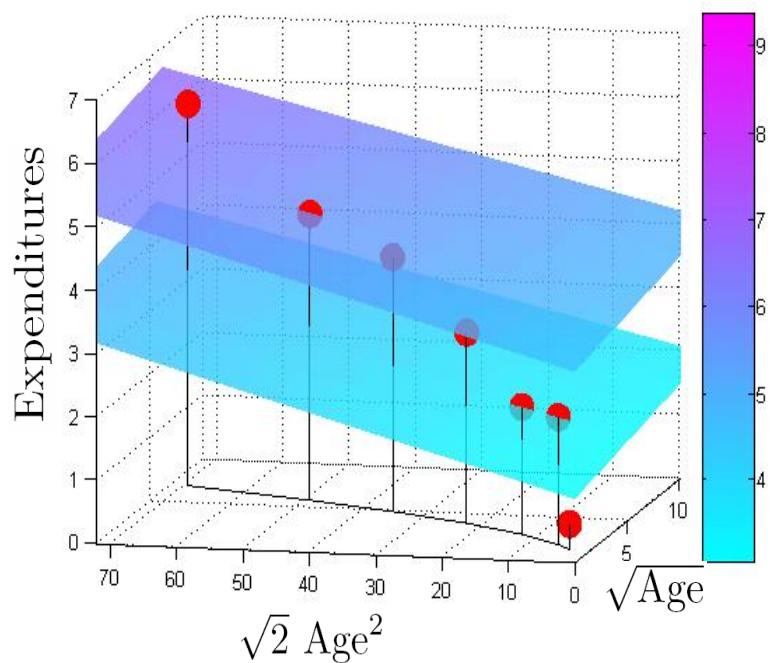
$$y_i - (w_1 x_{i1}) - b \leq \epsilon + \xi_i \quad \text{\textbackslash\textbackslash\textbackslash}$$

$$(w_1 x_{i1}) + b - y_i \leq \epsilon + \xi_i^* \quad \text{\textbar\textbar\textbar}$$

$$\xi_i, \xi_i^* \geq 0 \quad i = 1, 2, \dots, n$$

Slide Credit: Georgi Nalbantov, Erasmus Research Institute of Management

Non-linear (Kernel) SVR



$$\min_{w_1, b, \xi_i, \xi_i^*} \frac{w_1^2 + w_2^2}{2} + C \sum_i (\xi_i + \xi_i^*)$$

Subject to:

$$y_i - (\mathbf{w}'\phi(x_{i1})) - b \leq \epsilon + \xi_i$$

$$(\mathbf{w}'\phi(x_{i1})) + b - y_i \leq \epsilon + \xi_i^*$$

$$\xi_i, \xi_i^* \geq 0 \quad i = 1, 2, \dots, n$$

$$y = \mathbf{w}'\Phi(x) + b$$

Slide Credit: Georgi Nalbantov, Erasmus Research Institute of Management

SVR: Derivation

Subject to:

$$\min_{\mathbf{w}, b, \xi_i, \xi_i^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*)$$
$$y_i - (\mathbf{w}' \phi(\mathbf{x}_i)) - b \leq \epsilon + \xi_i$$
$$(\mathbf{w}' \phi(\mathbf{x}_i)) + b - y_i \leq \epsilon + \xi_i^*$$
$$\xi_i, \xi_i^* \geq 0 \quad i = 1, 2, \dots, n$$

$$L := \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*) - \sum_i (\eta_i \xi_i + \eta_i^* \xi_i^*)$$
$$- \sum_i \alpha_i (\epsilon + \xi_i - y_i + \mathbf{w}' \phi(\mathbf{x}_i) + b) - \sum_i \alpha_i^* (\epsilon + \xi_i^* + y_i - \mathbf{w}' \phi(\mathbf{x}_i) - b)$$

min with respect to $\mathbf{w}, b, \xi_i, \xi_i^*$
max with respect to $\alpha_i, \alpha_i^*, \eta_i, \eta_i^*$

Slide Credit: Georgi Nalbantov, Erasmus Research Institute of Management

SVR: Summary

- **Strengths of SVR:**

- No local minima
- Scales relatively well to high-dimensional data
- Trade-off between classifier complexity and error can be controlled explicitly via C and ϵ
- Overfitting is avoided (for any fixed C and ϵ)
- The “curse of dimensionality” is avoided through kernel functions

- **Weaknesses of SVR:**

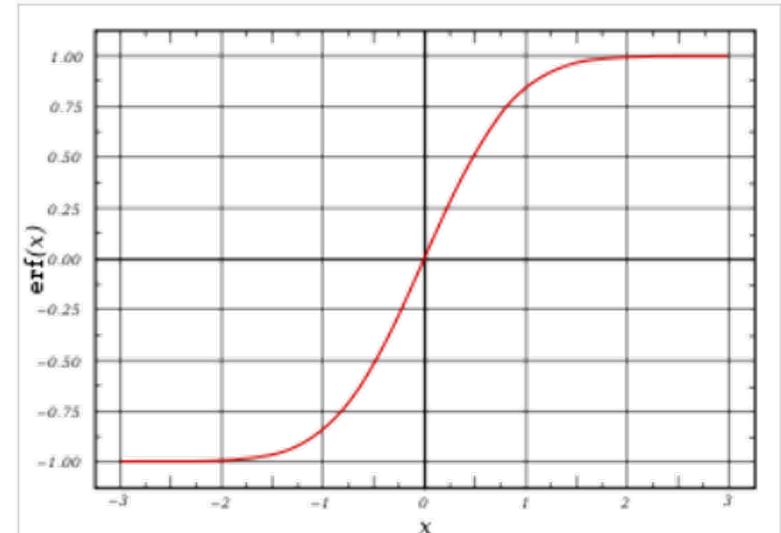
- What is the best trade-off parameter C and best ϵ ?
- What is a good transformation of the original space?

Regression Methods

- Linear Least-Squares Regression
 - Partial Least-Squares
 - Total Least-Squares
 - Ridge Regression, LASSO
- Kernel Regression
- k-NN Regression
- Regression Trees
- Support Vector Regression
- Logistic Regression

Logistic Regression

- To predict an outcome variable that is categorical from one or more categorical or continuous predictor variables.
- Used because having a categorical outcome variable violates the assumption of linearity in normal regression.
- Let X be the data instance, and Y be the class label: Learn $P(Y|X)$ directly
 - Let $W = (W_1, W_2, \dots, W_n)$, $X = (X_1, X_2, \dots, X_n)$, $W \cdot X$ is the dot product
 - Sigmoid function:
$$P(Y = 1 | X) = \frac{1}{1 + e^{-W \cdot X}}$$



Logistic Regression

- Generative or Discriminative?

Logistic Regression

- Generative classifier, e.g., Naïve Bayes:
 - Assume some functional form for $\mathbf{P}(\mathbf{X}|\mathbf{Y}), \mathbf{P}(\mathbf{Y})$
 - Estimate parameters of $\mathbf{P}(\mathbf{X}|\mathbf{Y}), \mathbf{P}(\mathbf{Y})$ directly from training data
 - Use Bayes rule to calculate $\mathbf{P}(\mathbf{Y}|\mathbf{X}=\mathbf{x})$
 - This is ‘generative’ model
 - Indirect computation of $\mathbf{P}(\mathbf{Y}|\mathbf{X})$ through Bayes rule
 - But, can generate a sample of the data, $P(X) = \sum_y P(y)P(X|y)$
- Discriminative classifier, e.g., Logistic Regression:
 - Assume some functional form for $\mathbf{P}(\mathbf{Y}|\mathbf{X})$
 - Estimate parameters of $\mathbf{P}(\mathbf{Y}|\mathbf{X})$ directly from training data
 - This is the ‘discriminative’ model
 - Directly learn $\mathbf{P}(\mathbf{Y}|\mathbf{X})$

Logistic Regression

- In logistic regression, we learn the conditional distribution $P(y|x)$
- Let $p_y(x;w)$ be our estimate of $P(y|x)$, where w is a vector of adjustable parameters.
- Assume there are two classes, $y = 0$ and $y = 1$ and

$$p_1(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}\mathbf{x}}} \quad p_0(\mathbf{x}; \mathbf{w}) = 1 - \frac{1}{1 + e^{-\mathbf{w}\mathbf{x}}}$$

- This is equivalent to $\log \frac{p_1(\mathbf{x}; \mathbf{w})}{p_0(\mathbf{x}; \mathbf{w})} = \mathbf{w}\mathbf{x}$
- That is, the log odds of class 1 is a linear function of x
- Q: How to find \mathbf{W} ?

Logistic Regression

- Conditional data likelihood - Probability of observed Y values in the training data, conditioned on corresponding X values.
- We choose parameters w that satisfy

$$\mathbf{w} = \arg \max_{\mathbf{w}} \prod_l P(y^l | \mathbf{x}^l, \mathbf{w})$$

- where
 - $\mathbf{w} = \langle w_0, w_1, \dots, w_n \rangle$ is the vector of parameters to be estimated,
 - y^l denotes the observed value of Y in the l th training example, and
 - \mathbf{x}^l denotes the observed value of \mathbf{X} in the l th training example

Logistic Regression

- Equivalently, we can work with log of conditional likelihood:

$$\mathbf{w} = \arg \max_{\mathbf{w}} \sum_l \ln P(y^l | \mathbf{x}^l, \mathbf{w})$$

- Conditional data log likelihood, $l(\mathbf{W})$, can be written as

$$l(\mathbf{w}) = \sum_l y^l \ln P(y^l = 1 | \mathbf{x}^l, \mathbf{w}) + (1 - y^l) \ln P(y^l = 0 | \mathbf{x}^l, \mathbf{w})$$

- Note here that Y can take only values 0 or 1, so only one of the two terms in the expression will be non-zero for any given y^l

Logistic Regression: Training

- We need to estimate:

$$\mathbf{w} = \arg \max_{\mathbf{w}} \sum_l \ln P(y^l | \mathbf{x}^l, \mathbf{w})$$

- Equivalently, we can minimize negative log likelihood
- This is convex – so, unique global minimum
- No closed-form solution though. Iterative method required.

Logistic Regression: Training

- Use gradient ascent (descent) for the maximization (min) problem
- The i th component of the vector gradient has the form

$$\frac{\partial}{\partial w_i} l(\mathbf{w}) = \sum_l x_i^l (y^l - \hat{P}(y^l = 1 | \mathbf{x}^l, \mathbf{w}))$$

Logistic Regression
prediction

- Beginning with initial weights, we repeatedly update the weights in the direction of the gradient, changing the i th weight according to

$$w_i \leftarrow w_i + \eta \sum_l x_i^l (y^l - \hat{P}(y^l = 1 | \mathbf{x}^l, \mathbf{w}))$$

Regularization in Logistic Regression

- Overfitting can arise especially when data has very high dimensions and is sparse.
- One approach -> modified “penalized log likelihood function,” which penalizes large values of \mathbf{w} , as before.

$$\mathbf{w} = \arg \max_{\mathbf{w}} \sum_l \ln P(y^l | \mathbf{x}^l, \mathbf{w}) - \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- Derivative then becomes:

$$\frac{\partial}{\partial w_i} l(\mathbf{w}) = \sum_l x_i^l (y^l - \hat{P}(y^l = 1 | \mathbf{x}^l, \mathbf{w})) - \lambda w_i$$

Logistic Regression: Summary

- In general, NB and LR make different assumptions
 - NB: Features independent given class -> assumption on $P(X|Y)$
 - LR: Functional form of $P(Y|X)$, no assumption on $P(X|Y)$
- LR is a linear classifier
 - decision rule is a hyperplane
- LR optimized by conditional likelihood
 - no closed-form solution
 - Concave (convex) -> global optimum with gradient ascent (descent)
- Extending logistic regression to multiple classes?
 - Use softmax for each class $k!$
$$p(y = k|x) = \frac{\exp(\theta_k^\top x)}{\sum_{i=1}^K \exp(\theta_i^\top x)}$$

Other Regression Methods

- Bayesian Regression
- Generalized Regression Neural Network

Readings

- PRML Bishop
 - Sec 3.1, Sec 4.3, Sec 7.1
- “Introduction to Machine Learning” by Ethem Alpaydin
 - Sec 4.6, 5.8, Chap 10, Chap 13
- Other resources
 - A nice tutorial on SVR: <https://alex.smola.org/papers/2004/SmoSch04.pdf>