

CS6510
Applied Machine Learning

Regression

1 Oct 2016

Vineeth N Balasubramanian

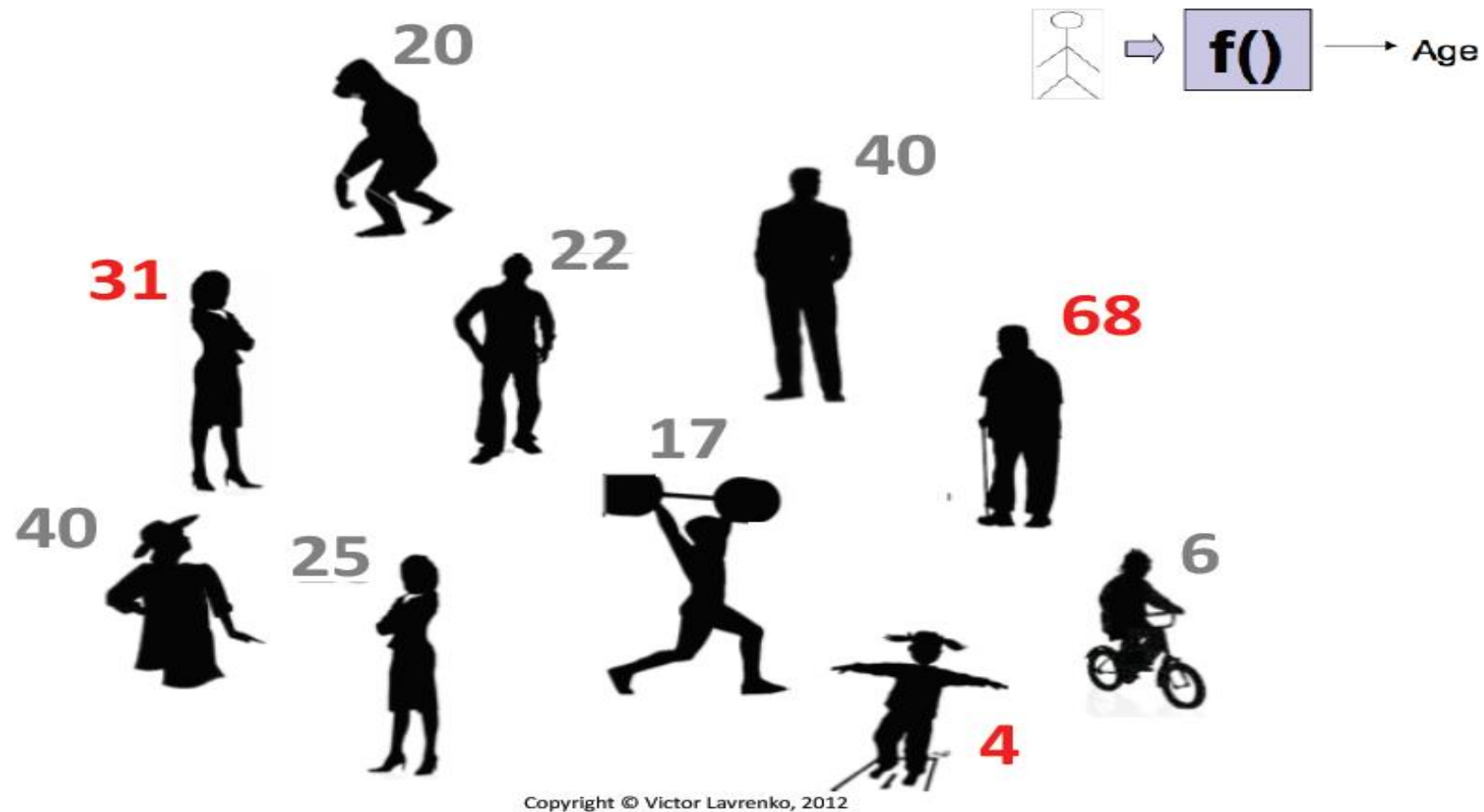


आई आई टी हैदराबाद
IIT Hyderabad

ML Problems

	<i>Supervised Learning</i>	<i>Unsupervised Learning</i>
<i>Discrete</i>	classification or categorization	clustering
<i>Continuous</i>	regression	dimensionality reduction

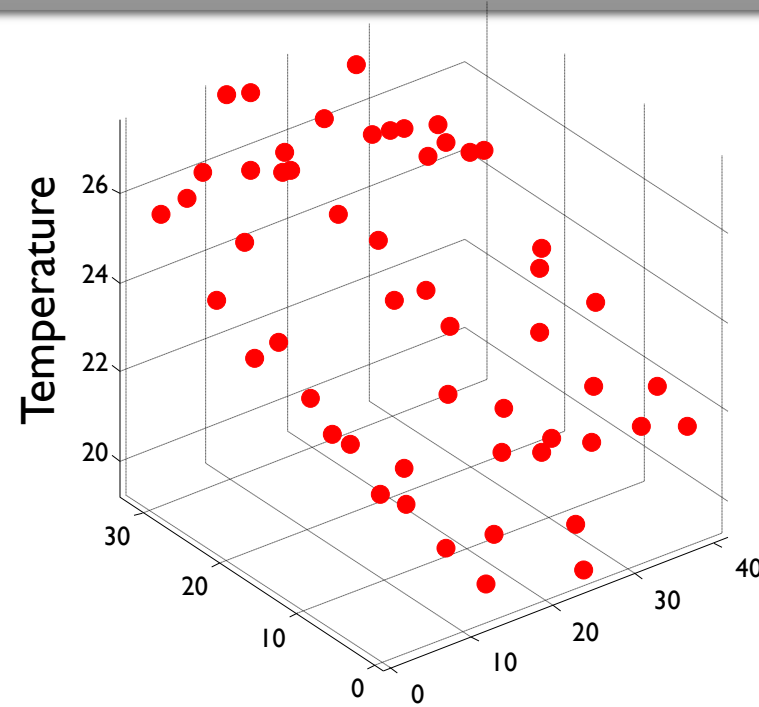
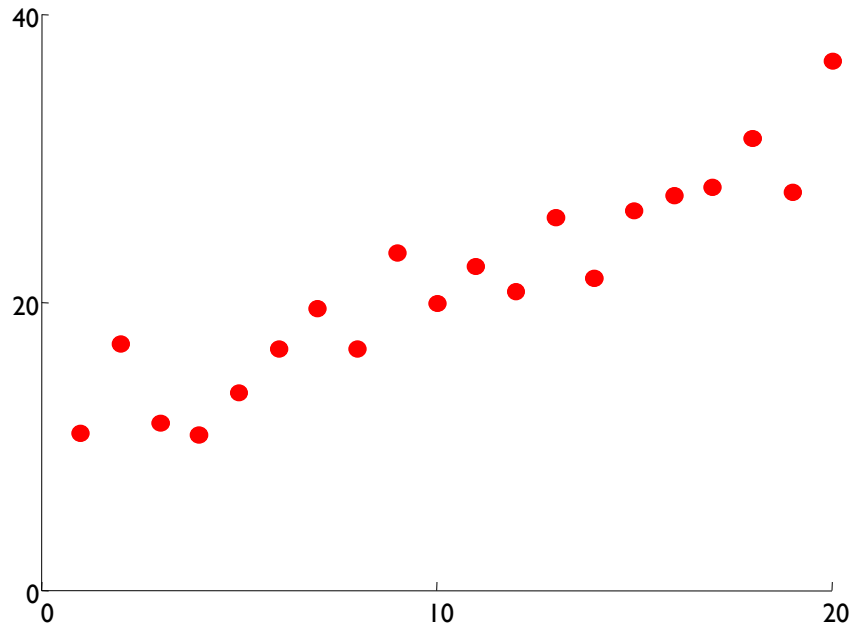
Regression (Supervised Learning)



Regression Methods

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

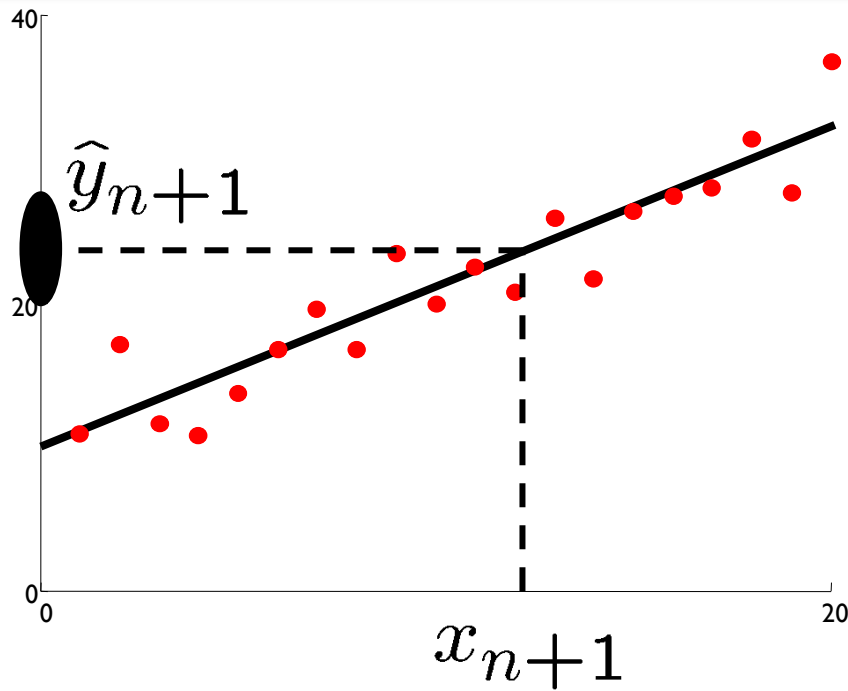
Linear Regression



Given examples $(x_i, y_i)_{i=1 \dots n}$

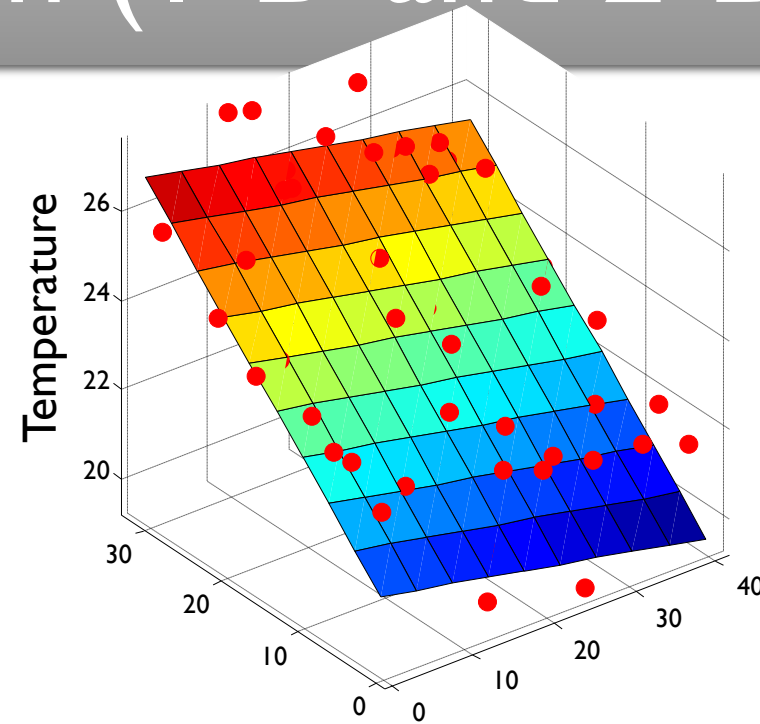
Predict y_{n+1} given a new point x_{n+1}

Linear Regression (1-D and 2-D)



Prediction

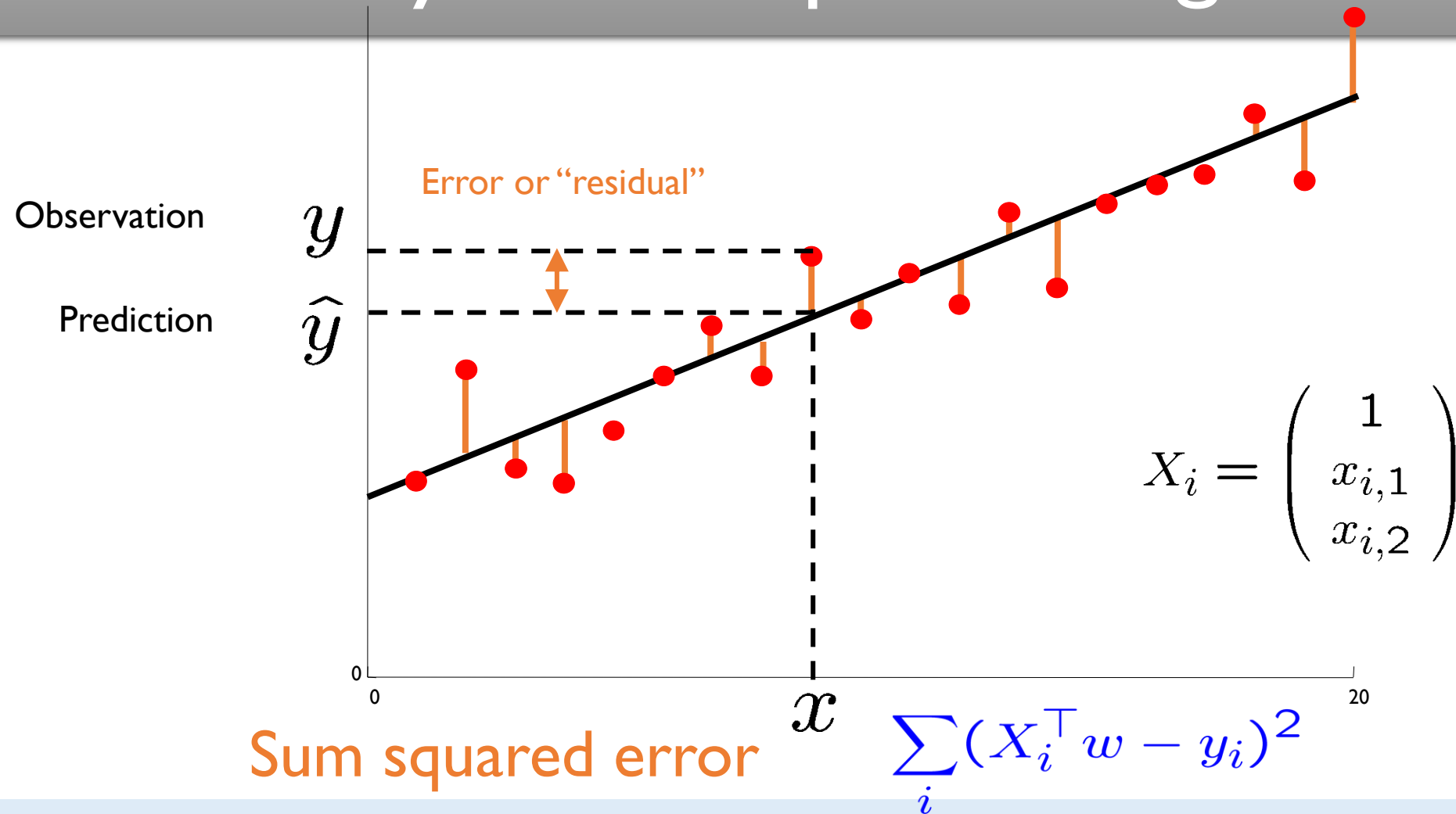
$$\hat{y}_i = w_0 + w_1 x_i$$



$$\text{Pre} = \begin{pmatrix} 1 & x_{i,1} & x_{i,2} \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix}$$

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

Ordinary Least-Squares Regression



Example

- Predict energy requirement

Wind Speed	People inside building	Energy requirement
100	2	5
50	42	25
45	31	22
60	35	18



Solving Least-Squares Regression

Sum squared error $E = \sum (X_i^\top w - y_i)^2$

$$E = \sum_i (X_i^\top w - y_i)^2$$

$$\begin{aligned}\frac{\partial E}{\partial w_j} &= \sum_i \frac{\partial}{\partial w_j} (X_i^\top w - y_i)^2 \\ &= \sum_i 2X_{i,j}(X_i^\top w - y_i)\end{aligned}$$

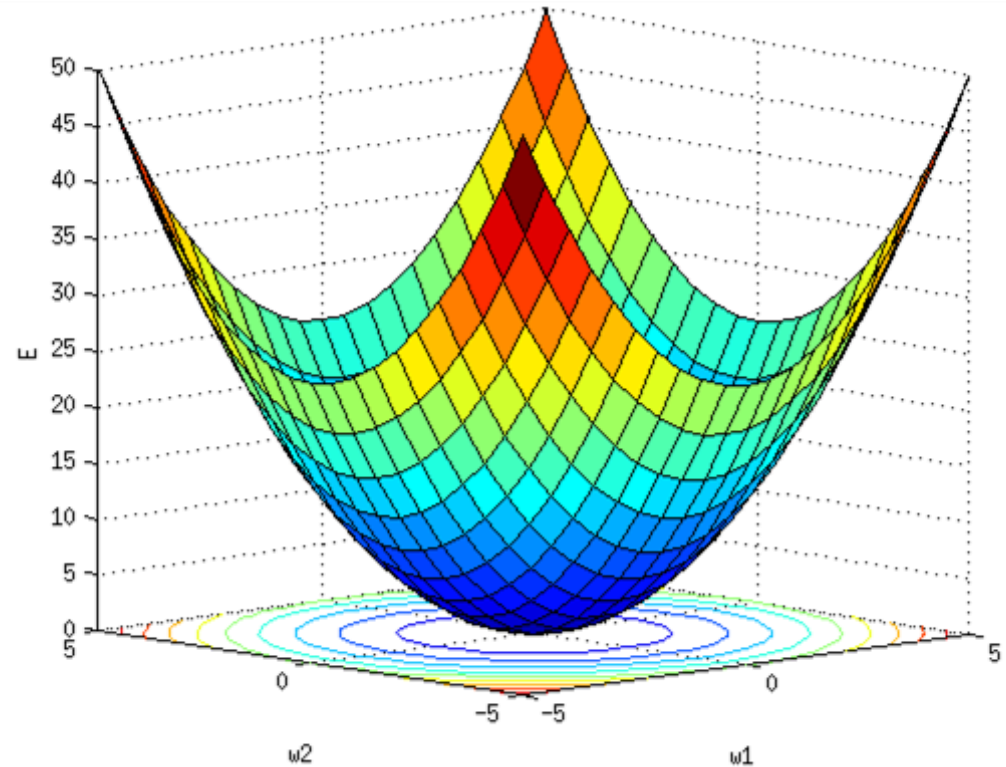
$$\frac{\partial E}{\partial w_j} = 0 \quad \Longleftrightarrow \quad \left(\sum_i X_{i,j} X_i^\top \right) w = \sum_i y_i X_{i,j} \quad \text{Linear equation}$$

$$\frac{\partial E}{\partial w} = 0 \quad \Longleftrightarrow \quad \underbrace{\left(\sum_i X_i X_i^\top \right)}_A w = \underbrace{\sum_i y_i X_i}_b \quad \begin{array}{l} \text{Linear system} \\ Aw = b \end{array}$$



Matrix View

- Sum of squared error:
 $E = (\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w})$
 - Why is this the same as what we saw?
 - What are the dimensions?
- What's the solution for this?
- Can we solve it by hand?
- Handling multiple outputs
 - Can be trivially extended:
 $\mathbf{y} = \mathbf{X}\mathbf{w}$ to $\mathbf{Y} = \mathbf{X}\mathbf{W}$



Linear Least Squares Regression

- With slight notation changes, we have linear least-squares solution to be:

$$\mathbf{X}'\mathbf{X}\mathbf{w} = \mathbf{X}'\mathbf{y}$$

- Assume $(\mathbf{X}'\mathbf{X})^{-1}$ exists, then

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

- Alternative Representation:

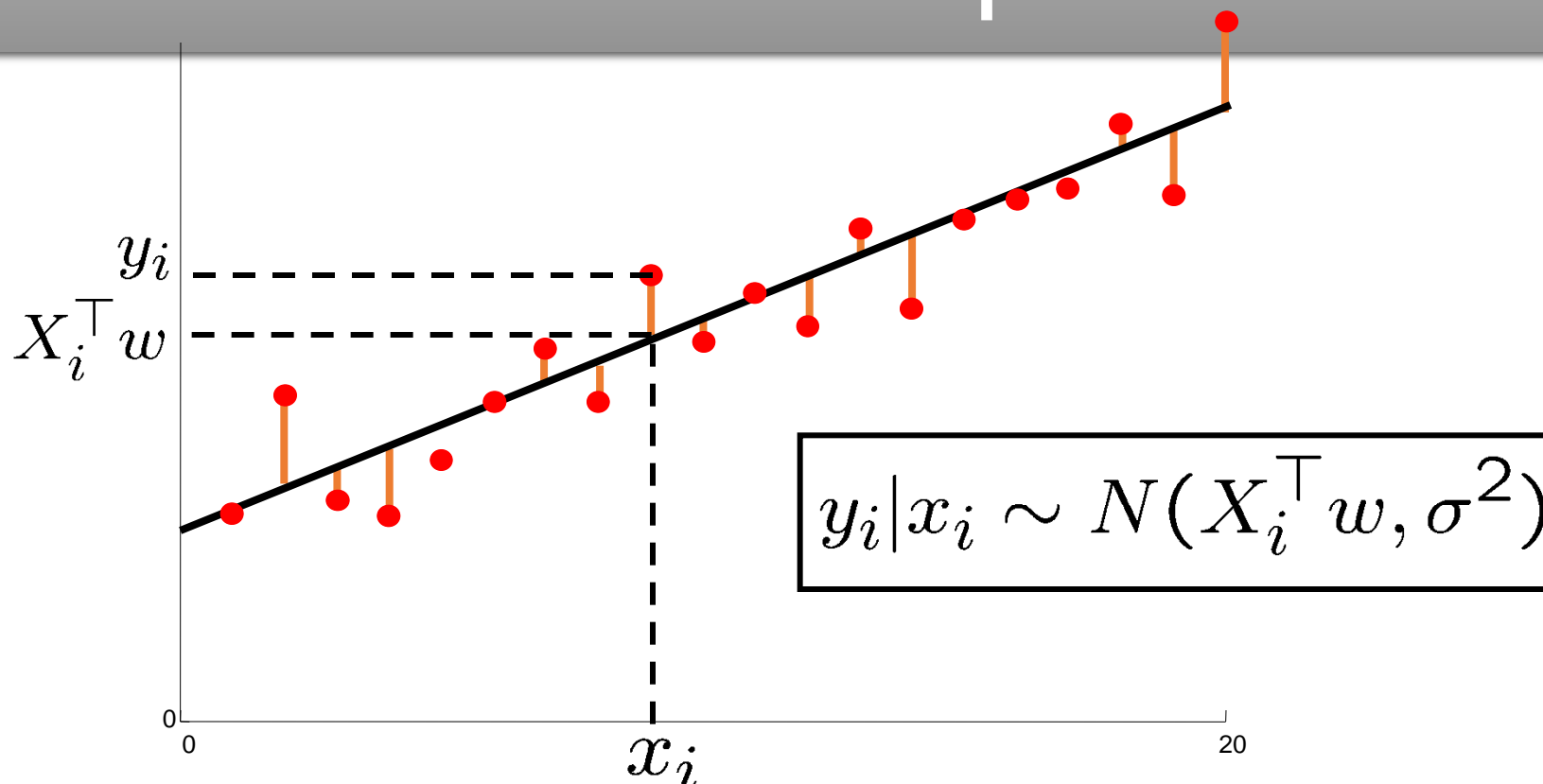
$$\begin{aligned}\mathbf{w} &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} = \mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} \\ &= \mathbf{X}'\left(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-2} \mathbf{X}'\mathbf{y}\right) = \mathbf{X}'\boldsymbol{\alpha}\end{aligned}$$

$$\text{where } \boldsymbol{\alpha} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-2} \mathbf{X}'\mathbf{y}, \mathbf{w} = \sum_{i=1}^{\ell} \alpha_i \mathbf{x}_i$$

Is this a good solution?

Computing inverse is not trivial; may not exist at all

Probabilistic Interpretation



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

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

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:

$$\sum_{j=1}^N (y_j - \sum_{i=0}^d w_i x_i)^2 + \lambda \sum_{i=1}^d w_i^2$$

- λ is the regularization parameter: bigger λ imposes more constraint

- **L₁ regularization (a.k.a. lasso):** Find \mathbf{w} which minimizes:

$$\sum_{j=1}^N (y_j - \sum_{i=0}^d w_i x_i)^2 + \lambda \sum_{i=1}^d |w_i|$$

Regularized Least-Squared Regression

- Understanding norms and regularization
 - L_p -norms:



$$p = \infty$$



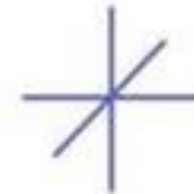
$$p = 2$$



$$p = 1$$



$$0 < p < 1$$



$$p = 0$$

Solving Ridge Regression

- Minimizing $\lambda \|\mathbf{w}\|^2 + \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$ 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

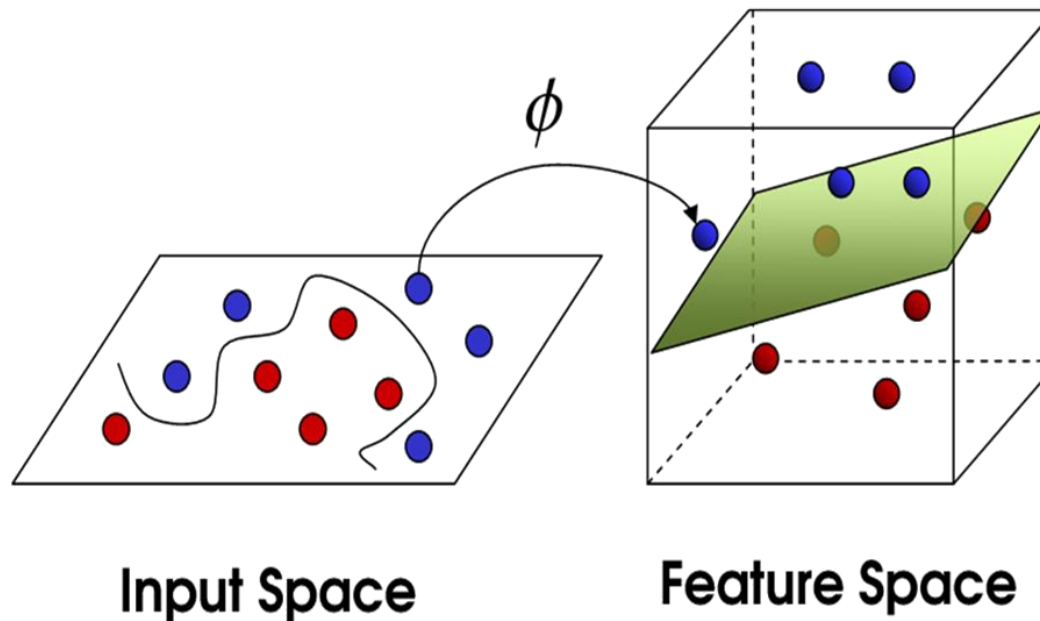
- **Total Least Squares:** Model error in output and input
 - Uses low-rank approximations to solve. E.g. SVD
- **Partial Least Squares:**
 - Seeks to address the correlation between predictor variables in its model
 - Also called “Projection to Latent Structures” (PLS)

Regression Methods

- Linear Least-Squares Regression
 - Partial Least-Squares
 - Total Least-Squares
 - Ridge Regression, LASSO
- Kernel Regression
- Logistic Regression
- k-NN Regression
- Regression Trees
- Support Vector 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 $O(d^3)$

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

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

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

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

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

$$\Rightarrow \boldsymbol{\alpha} = (\mathbf{G} + \lambda\mathbf{I}_\ell)^{-1} \mathbf{y} \text{ where } \mathbf{G} = \mathbf{X}\mathbf{X}'$$

Do you spot anything interesting?

Solving is $O(n^3)$

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

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

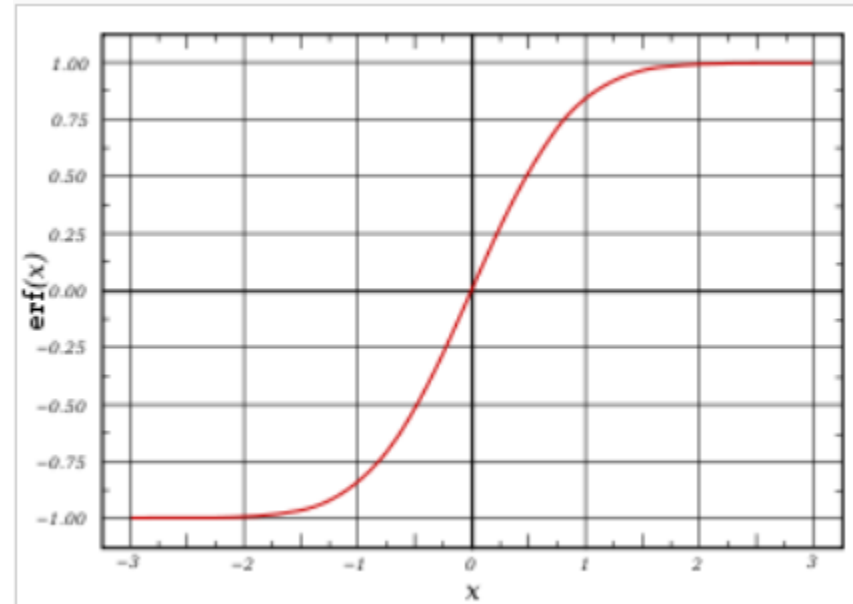
Regression Methods

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

Logistic Regression

- Let X be the data instance, and Y be the class label: Learn $P(Y|X)$ directly
 - Let $W = (W_1, W_2, \dots, W_n)$, $X = (X_1, X_2, \dots, X_n)$, $\mathbf{W}\mathbf{X}$ is the dot product
 - Sigmoid function:

$$P(Y = 1 | \mathbf{X}) = \frac{1}{1 + e^{-\mathbf{W}\mathbf{X}}}$$



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

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

- 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 **W**?

Logistic Regression

- The conditional data likelihood is the probability of the observed Y values in the training data, conditioned on their corresponding X values. We choose parameters \mathbf{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

- Learns the Conditional Probability Distribution $P(y|x)$
- Local Search.
 - Begins with initial weight vector.
 - Modifies it iteratively to maximize an objective function.
 - The objective function is the conditional log likelihood of the data – so the algorithm seeks the probability distribution $P(y|x)$ that is most likely given the data.

Logistic Regression

- In general, NB and LR make different assumptions
 - NB: Features independent given class \rightarrow 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 \rightarrow global optimum with gradient ascent

Regression Methods

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

k-NN Regression

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

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

Regression Methods

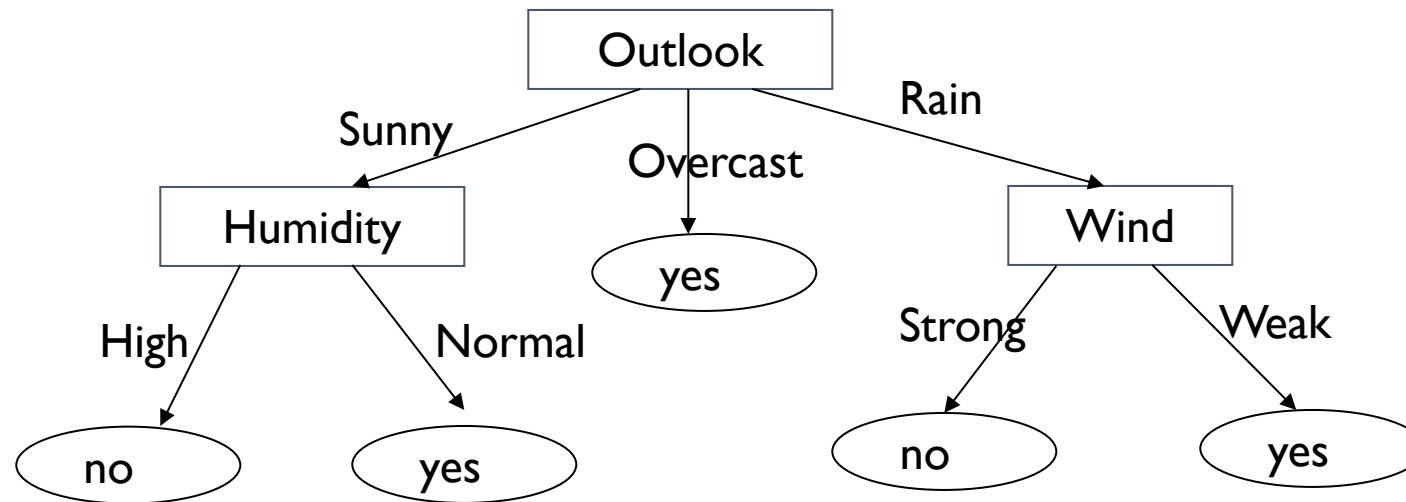
- Linear Least-Squares Regression
 - Partial Least-Squares
 - Total Least-Squares
 - Ridge Regression, LASSO
- Kernel Regression
- Logistic Regression
- k-NN Regression
- Regression Trees
- Support Vector Regression

Recall: Example

PlayTennis: training examples

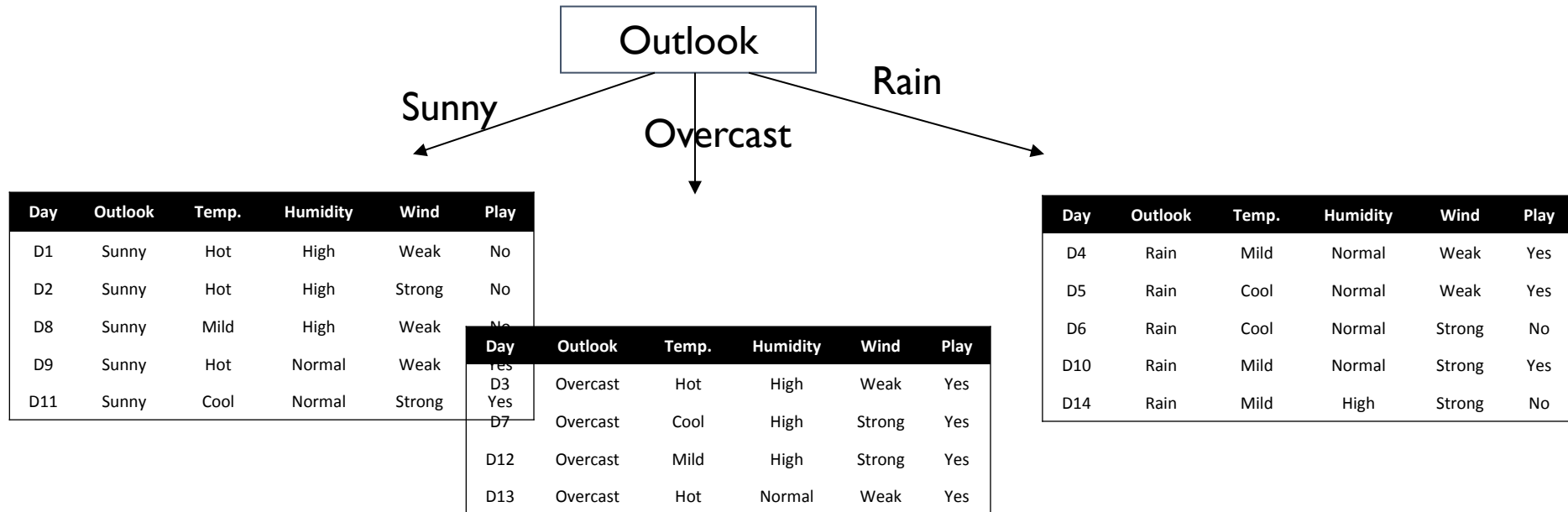
Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Recall: Example



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

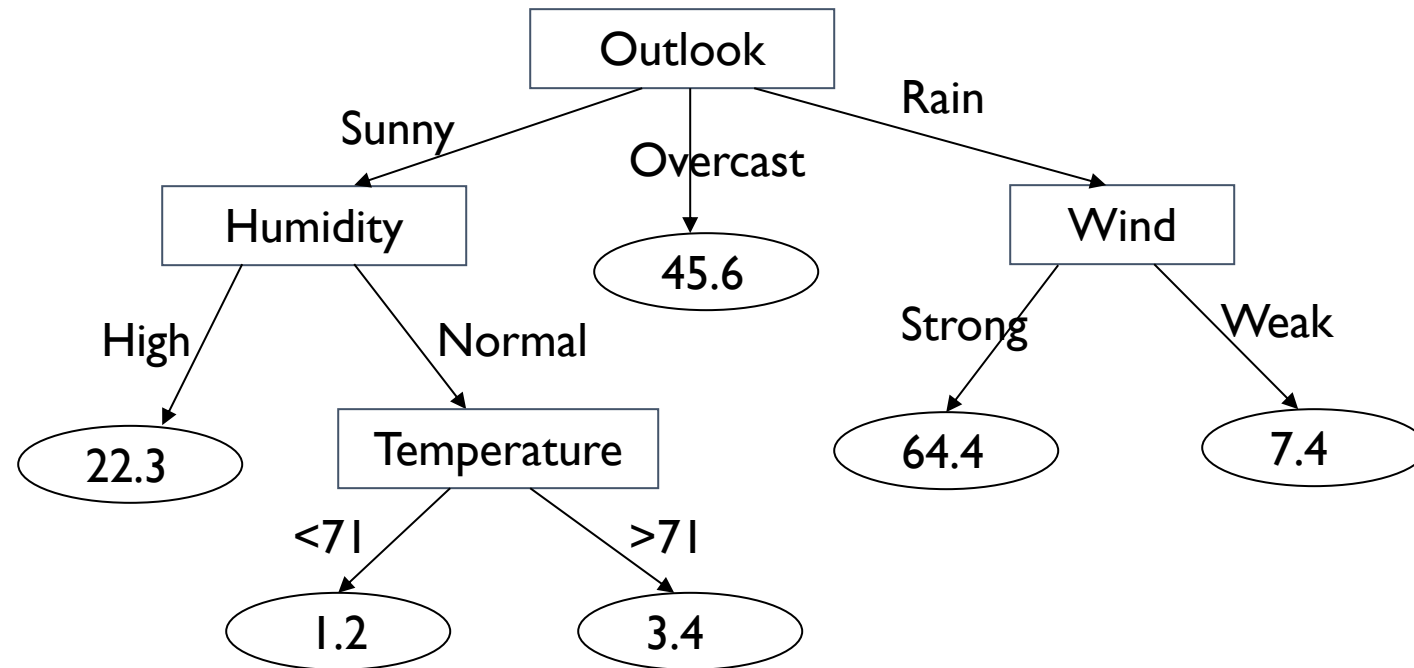
$$\Delta I(LS, A) = I(LS) - \sum_a \frac{|LS_a|}{|LS|} I(LS_a)$$

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

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

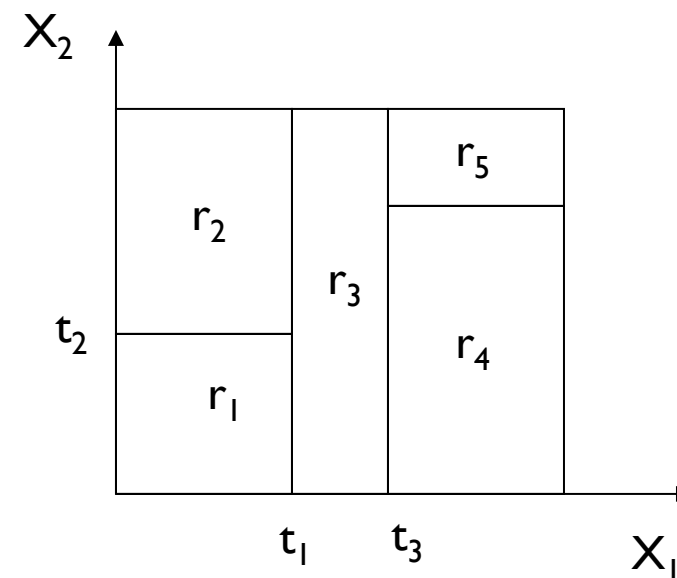
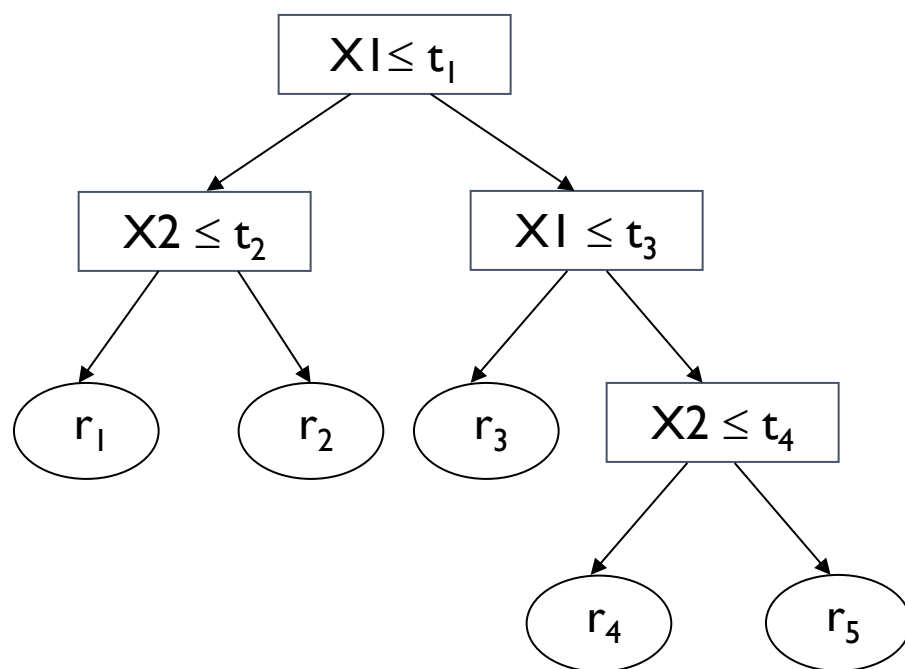
Regression Trees

- 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:

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

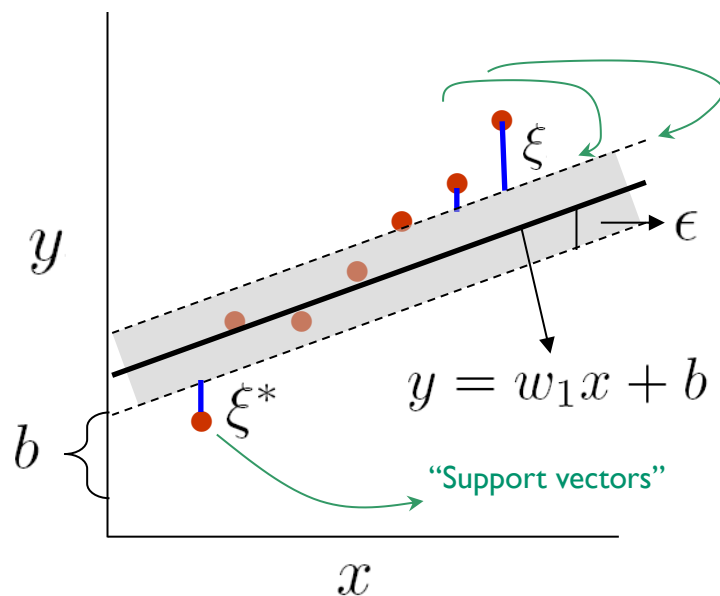
Regression Tree Pruning

- Exactly the same algorithms apply: pre-pruning and post-pruning.
- In post-pruning, the tree that minimizes the squared error on VS is selected.
- In practice, pruning is more important in regression because full trees are much more complex (often all objects have a different output values and hence the full tree has as many leaves as there are objects in the learning sample)

Regression Methods

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

Support Vector Regression



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

Complexity

Sum of errors

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

Subject to:

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

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

$$\xi_i, \xi_i^* \geq 0$$

Other Regression Methods

- Bayesian Regression
- Generalized Regression Neural Network