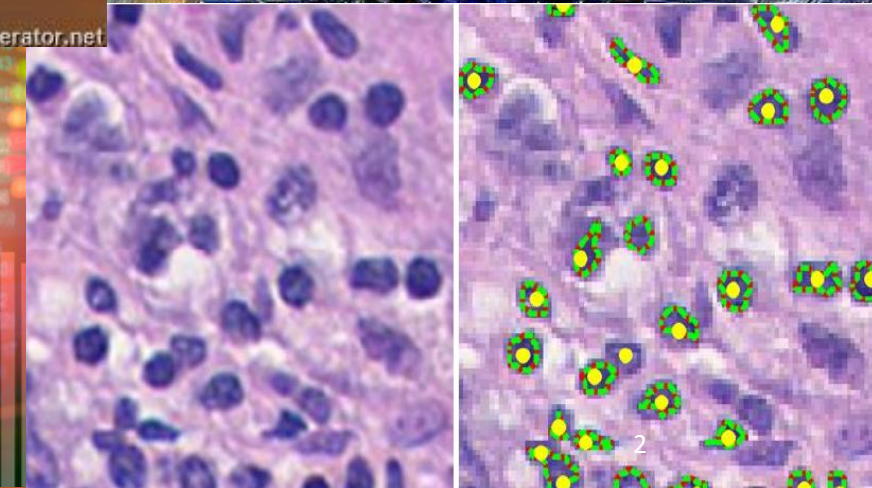
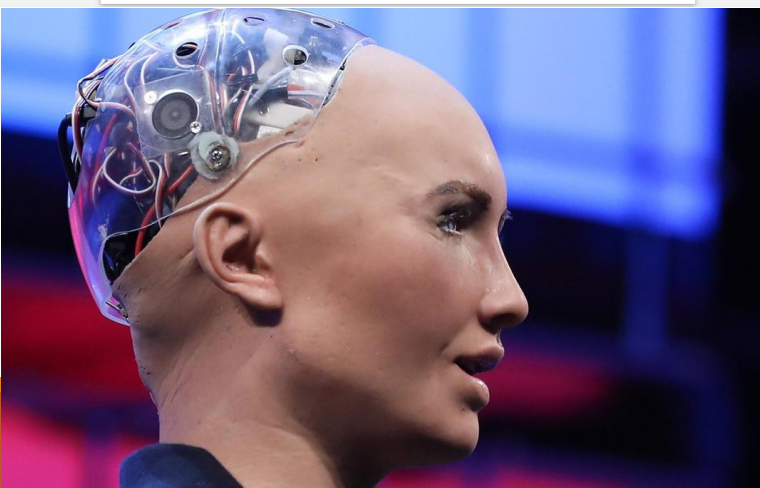
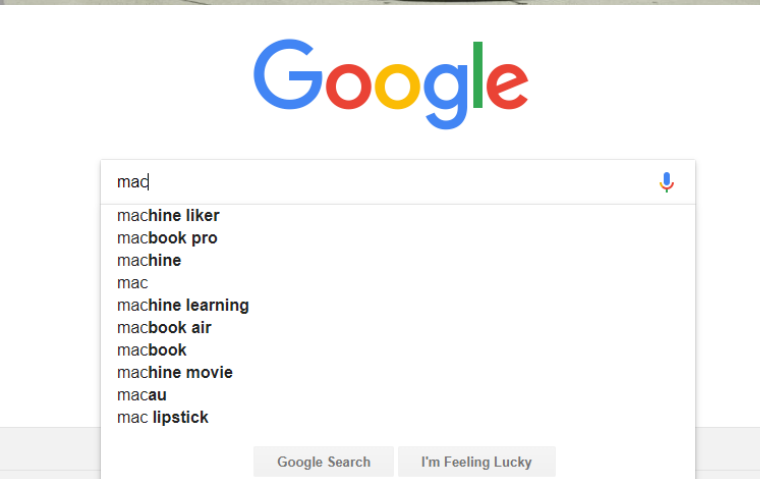
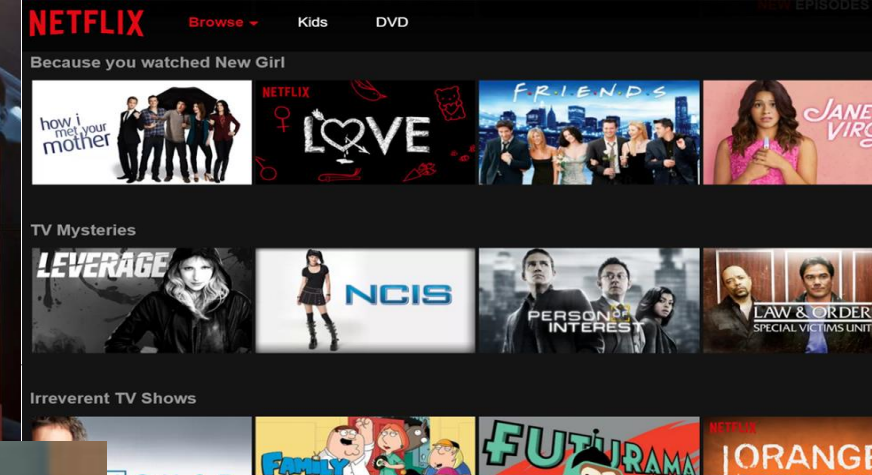


ML-101:

An Introduction to Machine Learning

BY SARTHAK CONSUL



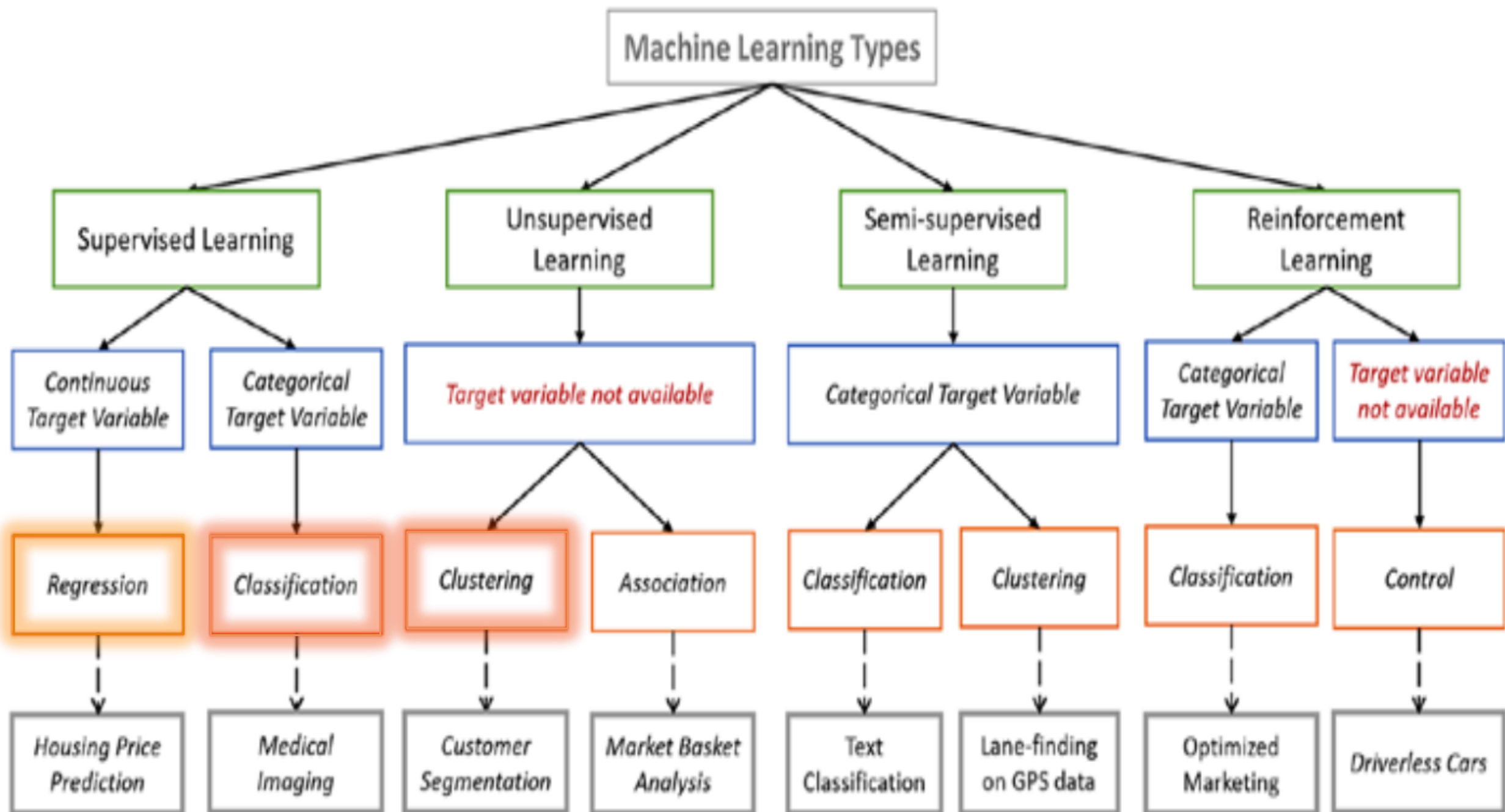
“Machine learning is glorified statistics”

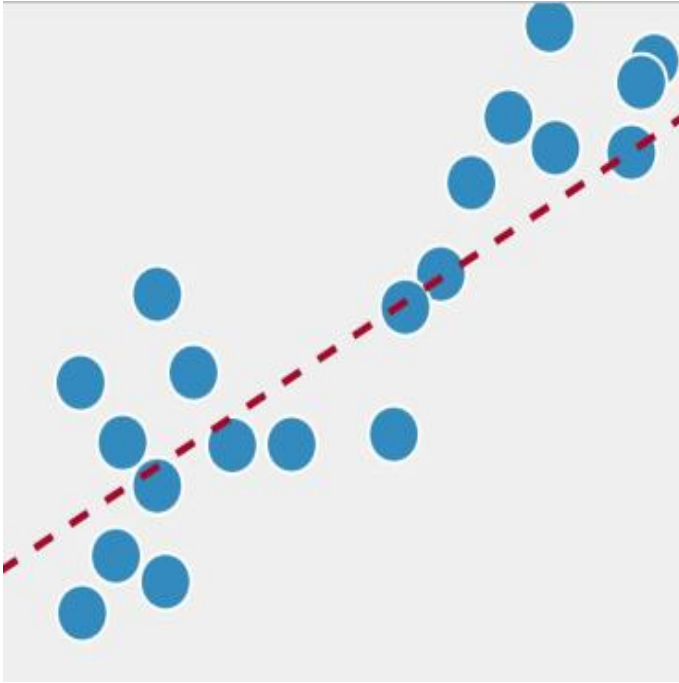
“Machine learning is statistics scaled up to big data”

“Machine learning is for Computer Science majors who couldn’t pass a Statistics course.” – A disgruntled statistician

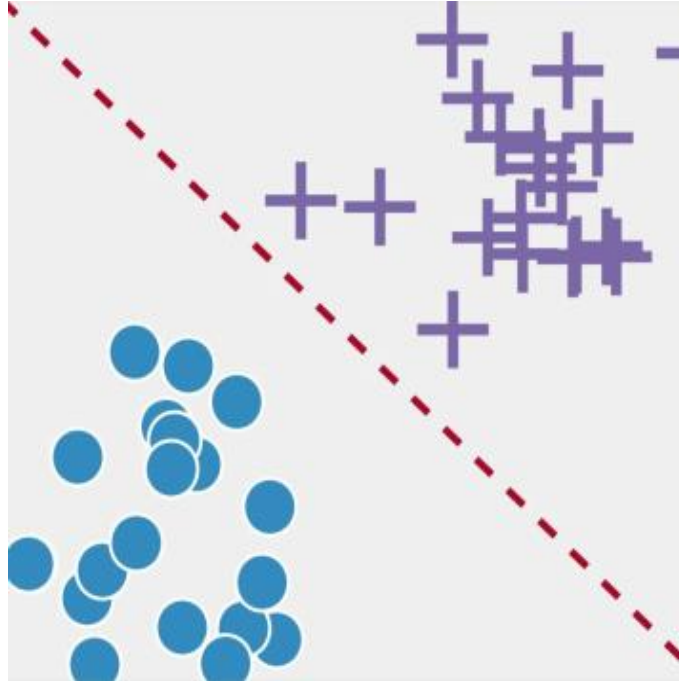
ML is an HYBRID field

- (Convex) Optimization
- Bayesian Statistics
- Lots and lots of linear algebra
- Sampling Theory
- Inspiration from biology, psychology

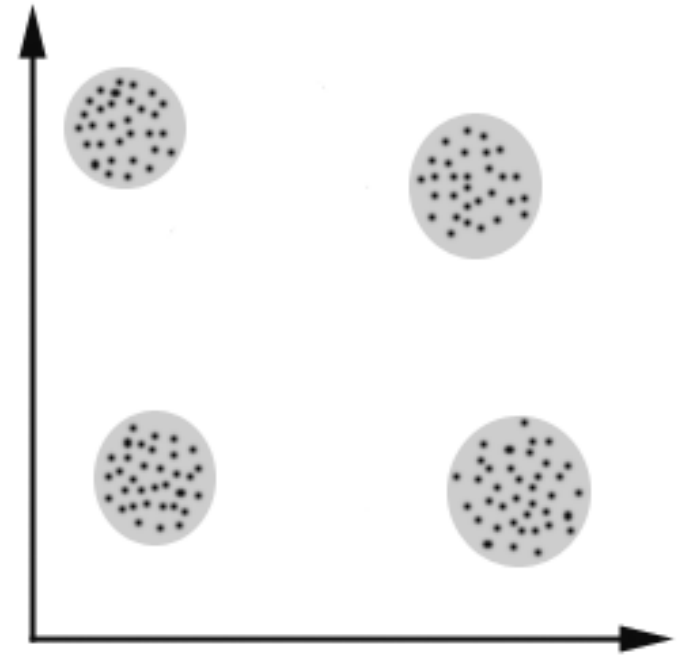




Regression



Classification

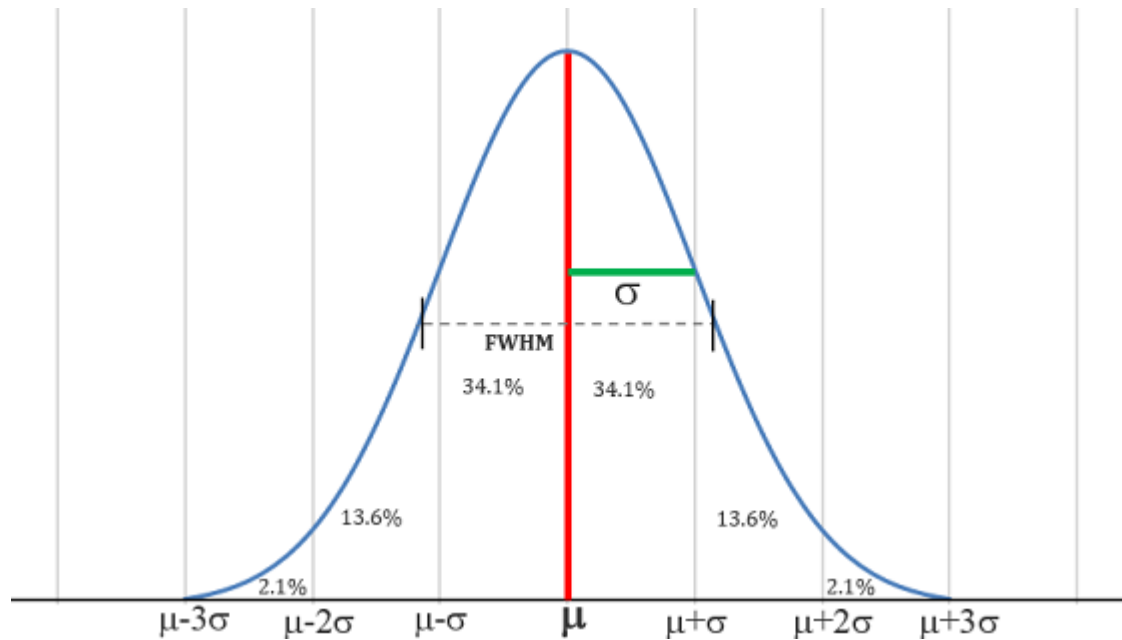


Clustering

Some Basic Statistics

❖ Probability Distributions

The probability of a real-valued variable x falling in the interval $(x, x + \delta x)$ is given by $p(x)\delta x$ for $\delta x \rightarrow 0$, then $p(x)$ is called the probability density over x



The Gaussian Distribution

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$$

Some Basic Statistics contd.

❖ Mean

$$\mathbb{E}[f] = \sum_x p(x) f(x)$$

$$\mathbb{E}[f] = \int p(x) f(x) \, dx.$$

❖ Variance

$$\text{var}[f] = \mathbb{E} [(f(x) - \mathbb{E}[f(x)])^2]$$

$$\text{var}[f] = \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2.$$

❖ Covariance

$$\begin{aligned} \text{cov}[x, y] &= \mathbb{E}_{x,y} [\{x - \mathbb{E}[x]\} \{y - \mathbb{E}[y]\}] \\ &= \mathbb{E}_{x,y} [xy] - \mathbb{E}[x]\mathbb{E}[y] \end{aligned}$$

❖ Standard Deviation

Pre-processing of Data

- ❖ Remove noise
- ❖ Filtering and Sampling
- ❖ Feature Scaling and Zero-Mean [Normalization]
- ❖ Building dictionaries
- ❖ Encoding
- ❖ Handling missing data
- ❖ PCA*
- ❖ Whitening

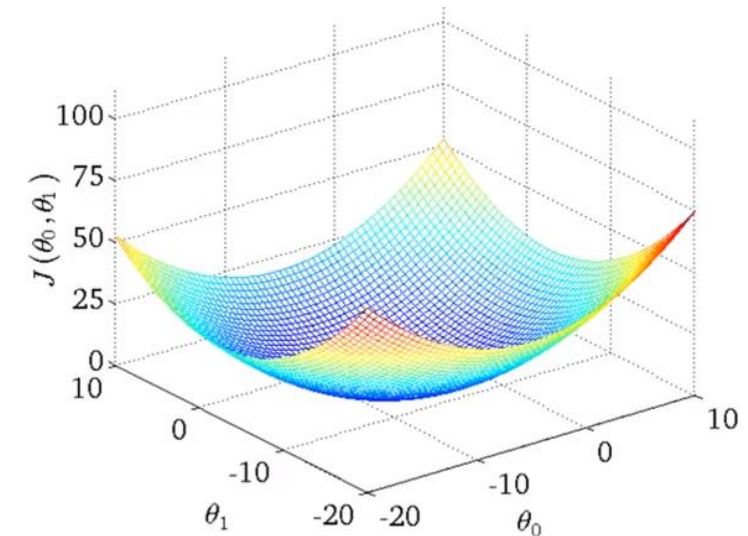
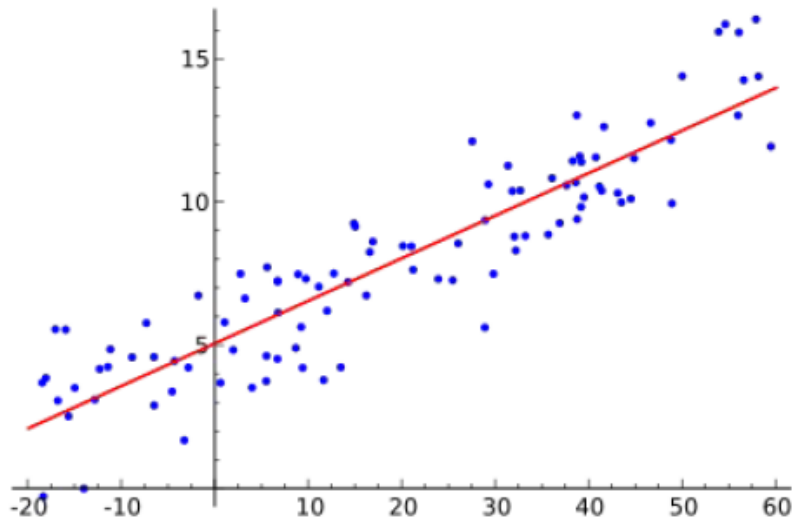
Linear Regression

❖ Hypothesis Function: $h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1.x_1 + \dots + \theta_n.x_n$

❖ Cost Function, $J(\theta)$ – Mean Squared Error

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

❖ Optimizing the cost Function via Gradient Descent



Gradient Descent

❖ Cost Function

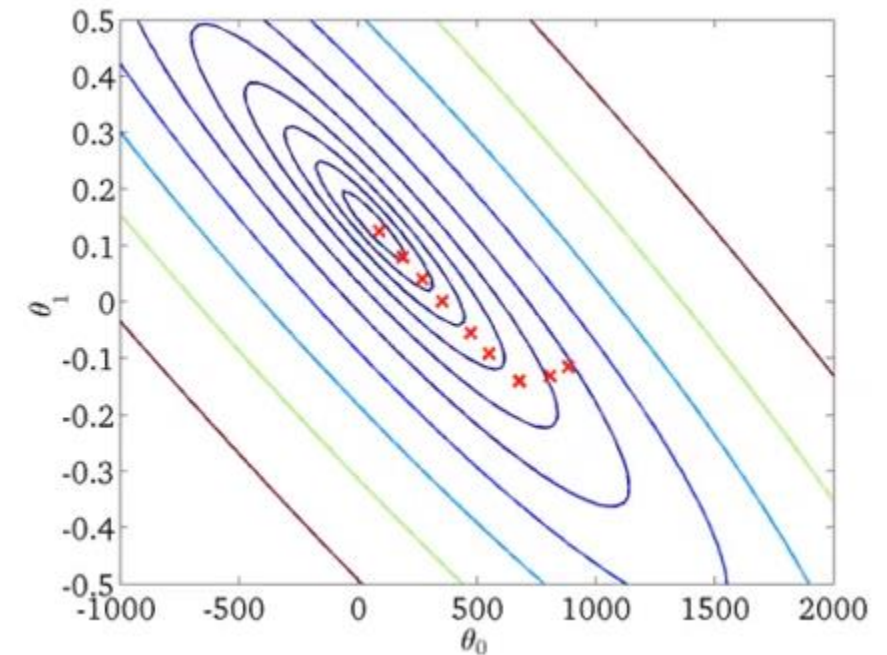
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

❖ Gradient Descent

Repeat till convergence{

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

}

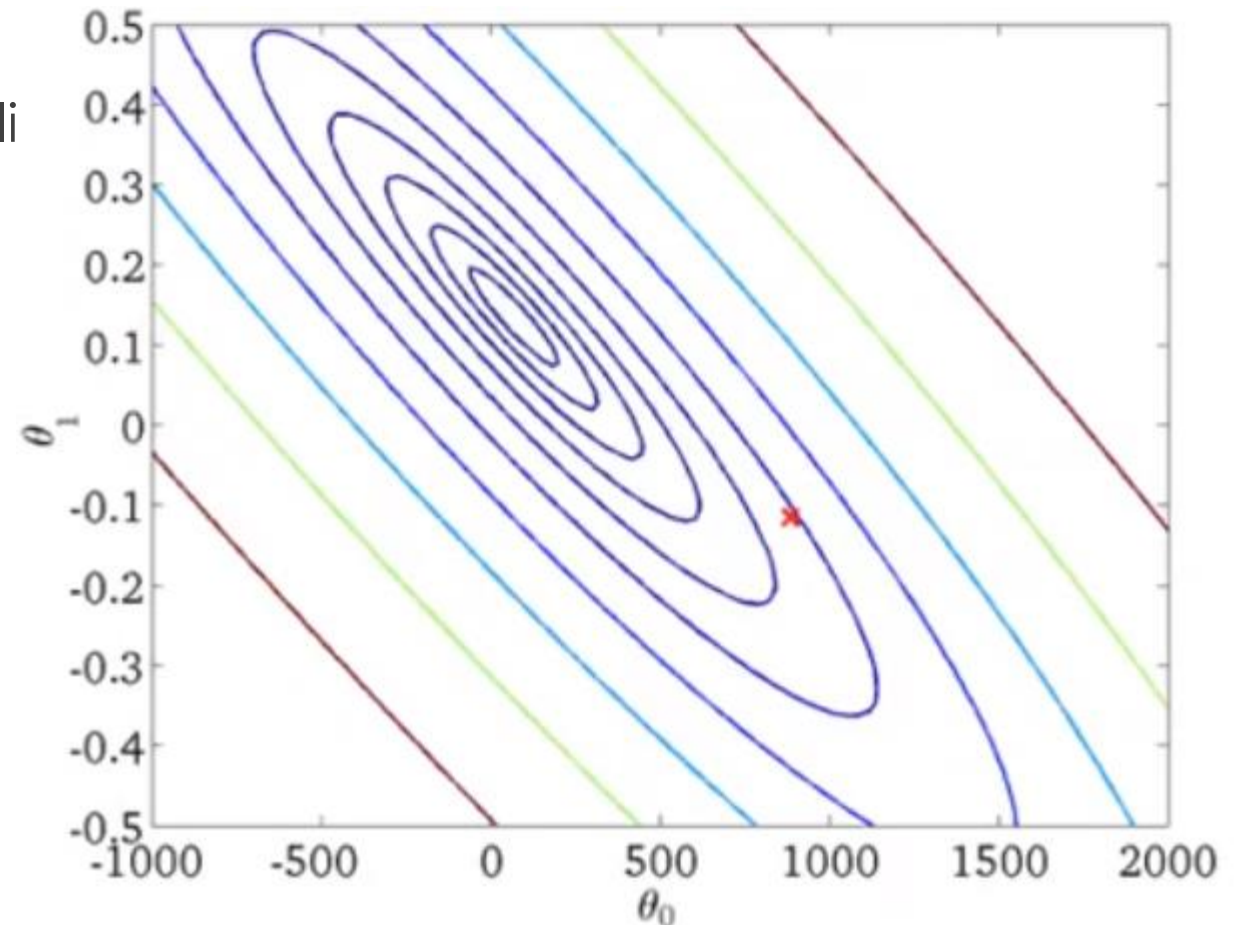


Gradient Descent contd.

- ❖ To improve convergence
 - ❖ Feature Scaling and Mean Normalization

$$x' = \frac{x - \bar{x}}{\sigma}$$

- ❖ Choosing the learning rate, α
- ❖ Batch Gradient Descent
- ❖ Stochastic Gradient Descent
- ❖ Mini-Batch Gradient Descent



Modifications to SGD

❖ SGD + Momentum

❖ Nesterov

❖ AdaGrad

❖ RMSProp

❖ Adam

...

```
dx = compute_gradient(x)
x += learning_rate * dx
```

```
dx = compute_gradient(x)
vx = rho * vx + dx
x += learning_rate * vx
```

Cross Validation

- ❖ While training model, test data should NOT be touched
- ❖ Limited Data
- ❖ Split Data into train, **validation**, (test)
- ❖ Holdout Method, K-fold CV, Leave-One-Out CV