Lecture - 10 (Logistic Regression)

Probabilistic :-

$$P(Y=y/X=x) = \frac{1}{1+\exp(y\langle\theta,x\rangle)}$$

$$y = \boxed{+1} \text{ or } \boxed{-1} = \frac{1}{1+\exp(\langle\theta,x\rangle)} \quad \frac{1}{1+\exp(-1\langle\theta,x\rangle)}$$

$\to$ $P(Y=1/x) + P(Y=-1/x=1) = 1$

$\to$ Decision boundary $\quad P(Y=1/x) = P(Y=-1/x) = 0.5$

$\to$ $P(Y=1/x) > 0.5 \Rightarrow$ greater confidence label 1 than -1

Prediction Confidence

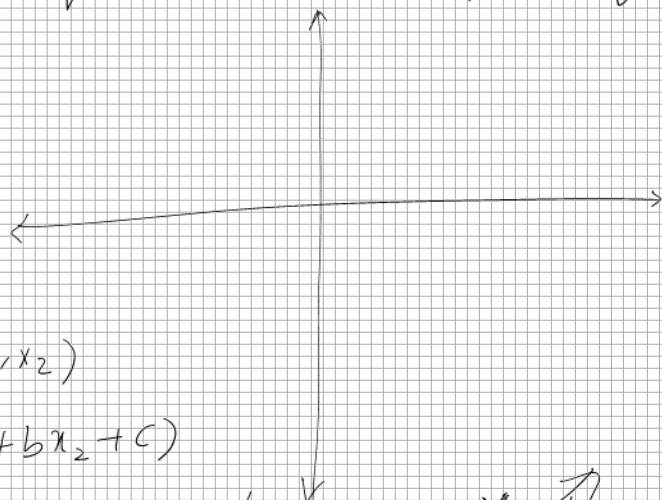| Task | Or |
|---|---|
| $\to$ predict the label associated with a featurevector $x$ | $\to$ predict grade $\text{Sign}\langle\theta,x\rangle$ |
| $\to$ Measure of confidence of that prediction | $\to$ $P(Y=y/X=x) = \frac{1}{1+\exp(y\langle\theta,x\rangle)}$ |

Properties of Linear classifiers

$\to$ ① Easy to train   ④ Linear classifiers - Excel at higher dimensions
$\hookrightarrow$ attraction computation load

$\to$ ② can predict labels very fast
( If $x$ is sparse, the inner product $\langle\theta,x\rangle$ is extremely fast )

$\to$ ③ Statistical theory of linear classifiers leading to Effective modeling strategies.

Bias Term $\quad \langle\theta,x\rangle + c$
$= x_1\theta_1 + x_2\theta_2 + \cdots + x_d\theta_d + \boxed{c}$

Bias term makes our model

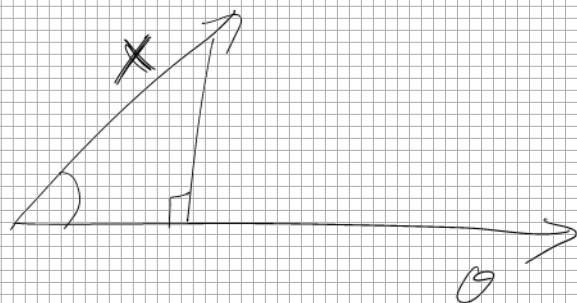→ Linear classifiers have the form $y = \text{sign}(\langle \theta, x \rangle)$

→ Given $x = (x_1, x_2)$

→ $\text{sign}(ax_1 + bx_2 + c)$

→ classification problem

decision boundary: don't
applied whether $\theta$ is normalized
or not.

→ Should the decision border pass
through (origin)?

$\hookrightarrow$ X or parameter
vector is zero.

Answer: Not Necessarily.
If we don't require that (by adding
our features to value 1) the classifier
becomes considerably more powerful.
i.e we can learn much more general
decision boundaries

$$X_1\theta_1 + X_2\theta_2 + \cdots + X_d\theta_d + c$$

without offset, the decision boundary has to
pass through Origin. i.e $\hookrightarrow \langle \theta, x \rangle = 0$

at origin $x=0, \theta=0 \Rightarrow \langle \theta, x \rangle = 0$

→ Decision boundary passing through origin restricts the model in a significant way.

$$n_1 \theta_1 + n_2 \theta_2 + \cdots + n_d \theta_d + C = \langle x, \theta \rangle + C$$

$$\langle x, \theta \rangle$$
↳ $\theta$ is $d+1$ dimensions.

$$(n_1, n_2 \cdots n_d, \underline{1})$$

---

→ # MLE

**Frequentists:** MLE uses pairs of feature vectors & labels, usually from historic data to estimate correct classifier or nature.

**Bayesian:** A single classifier cannot represent the "Truth". Estimate the posterior probability that each classifier is correct & use them all.

**Justification for using MLE:**
→ It converges to the optimal solution in the limit of large data (consistency).
→ Convergence occurs at the fastest possible rate of convergence (statistical efficiency).

**Gradient Descent:** is scalable to large data as long for Logistic. it is sparse, not so much dependent on the dimensionality.
dependent $\langle \theta, x \rangle$ = easy to calculate

$$\theta_j \longleftarrow \theta_j - \alpha \frac{\partial \sum_{i=1}^{n} \log \left(1 + \exp\left(y^{(i)} \langle \theta, x^{(i)} \rangle\right)\right)}{\partial \theta_j}$$

when $x^{(i)}$ are sparse, partial derivative can be made particularly fast.

$\rightarrow$ Stochastic Gradient converges faster than Gradient descent.