

Maximum Likelihood Estimation

and

Generalized Linear Models

Prepared by: Joseph Bakarji

Probabilistic Interpretation of Linear Regression

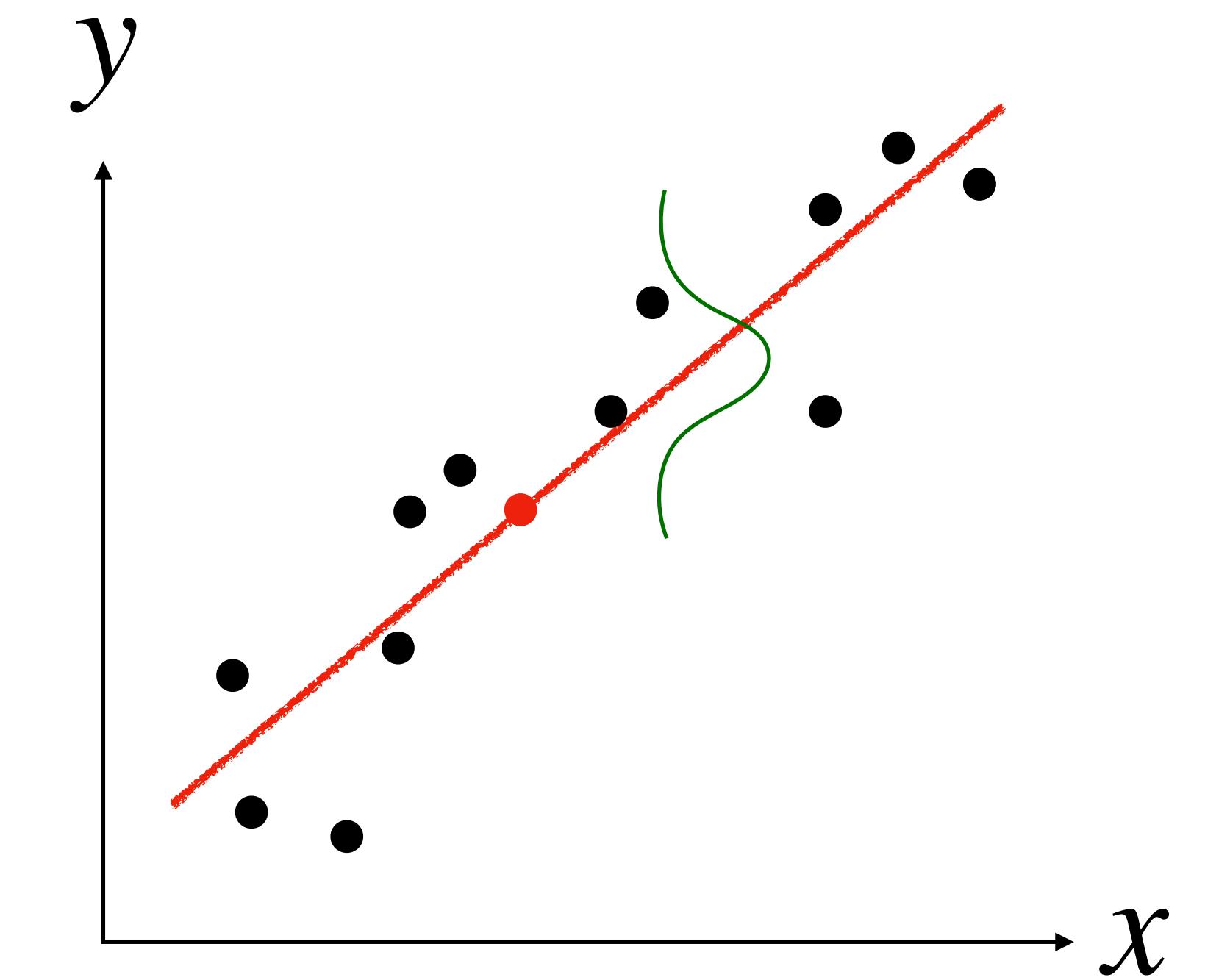
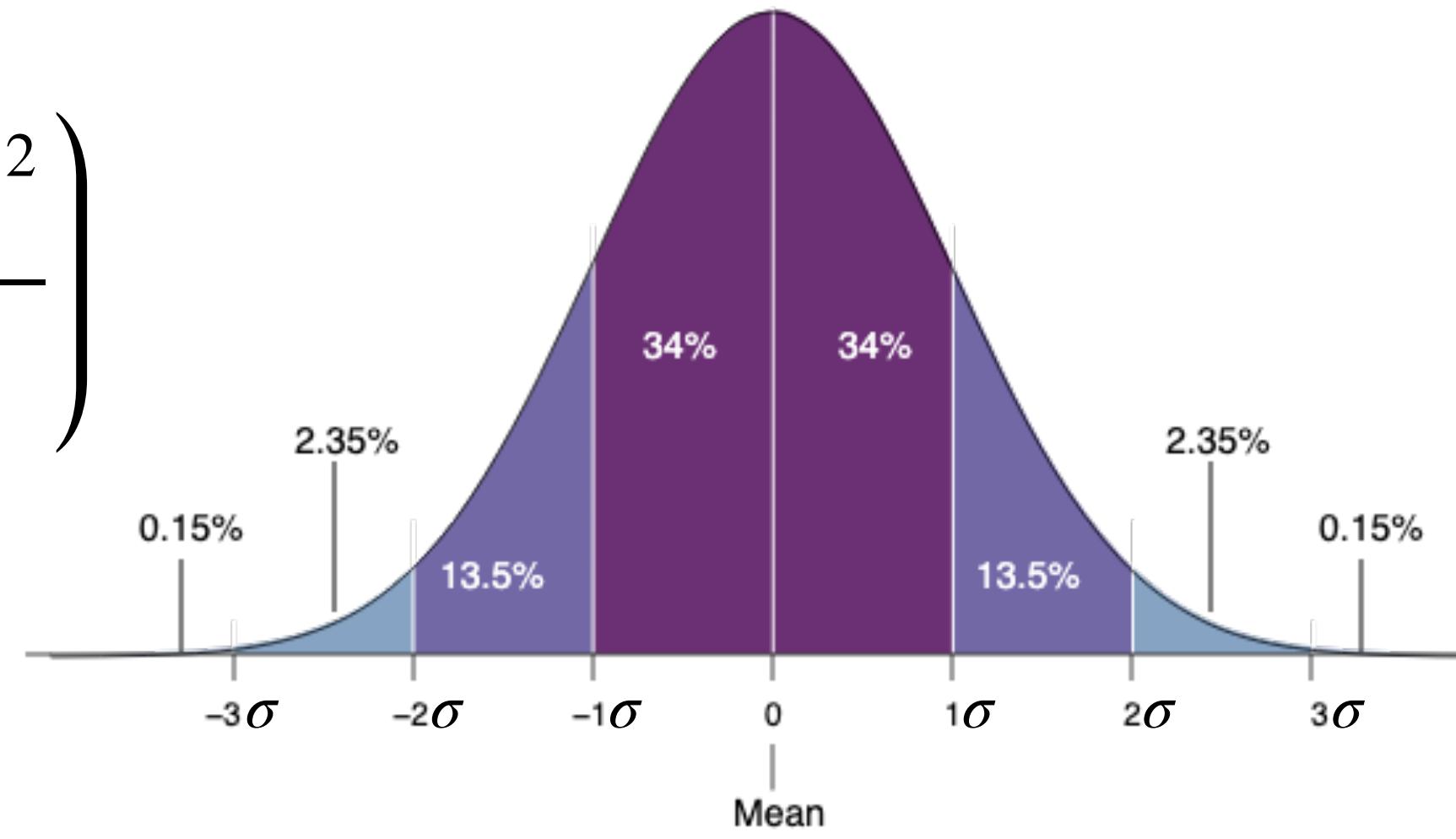
Assume **noise** is normally distributed around model

$$y^{(i)} = \theta^\top x^{(i)} + \varepsilon^{(i)}$$

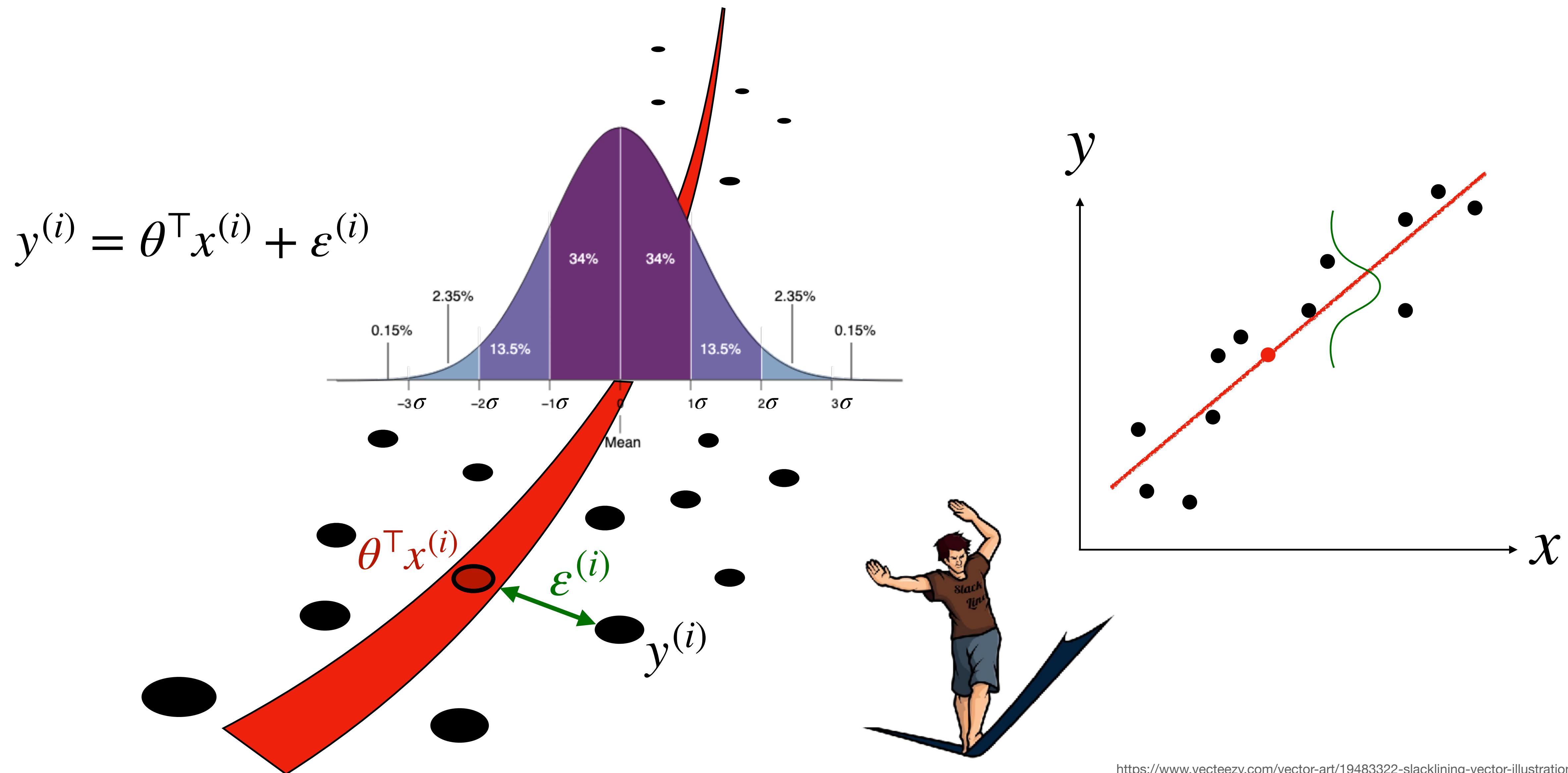
Normally distributed

$$\mathcal{N}(0, \sigma^2)$$

$$p(\varepsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\varepsilon^{(i)})^2}{2\sigma^2}\right)$$



Probabilistic Interpretation of Linear Regression

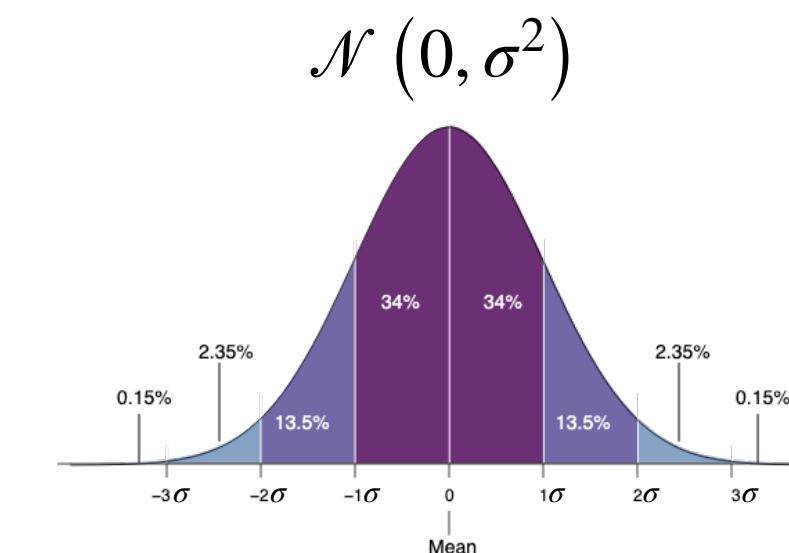


Probabilistic Interpretation of Linear Regression

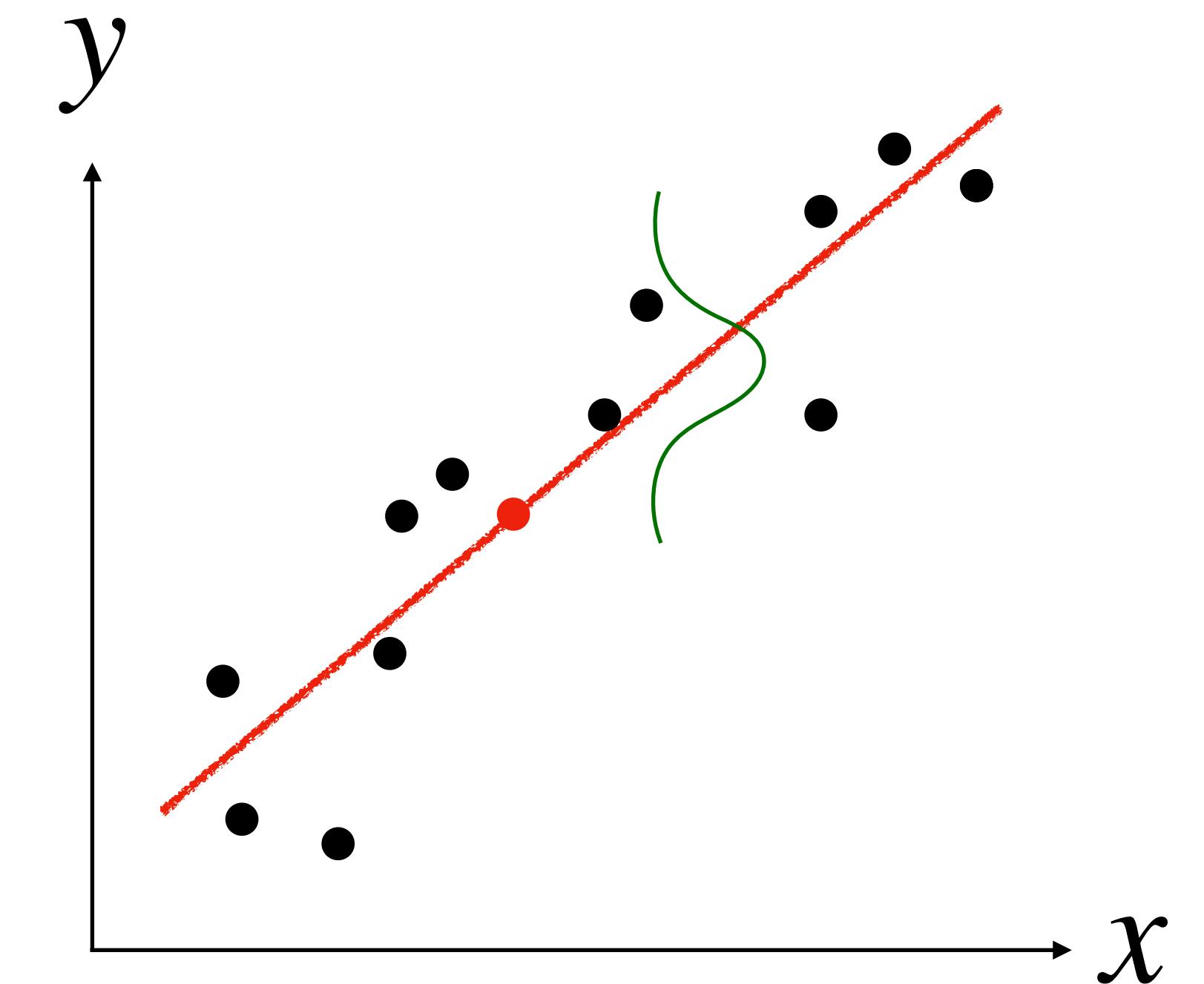
Assume noise is normally distributed around model

$$p(\varepsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\varepsilon^{(i)})^2}{2\sigma^2}\right)$$

$$y^{(i)} = \theta^\top x^{(i)} + \varepsilon^{(i)}$$



$$p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^\top x^{(i)})^2}{2\sigma^2}\right)$$



Likelihood of output given input

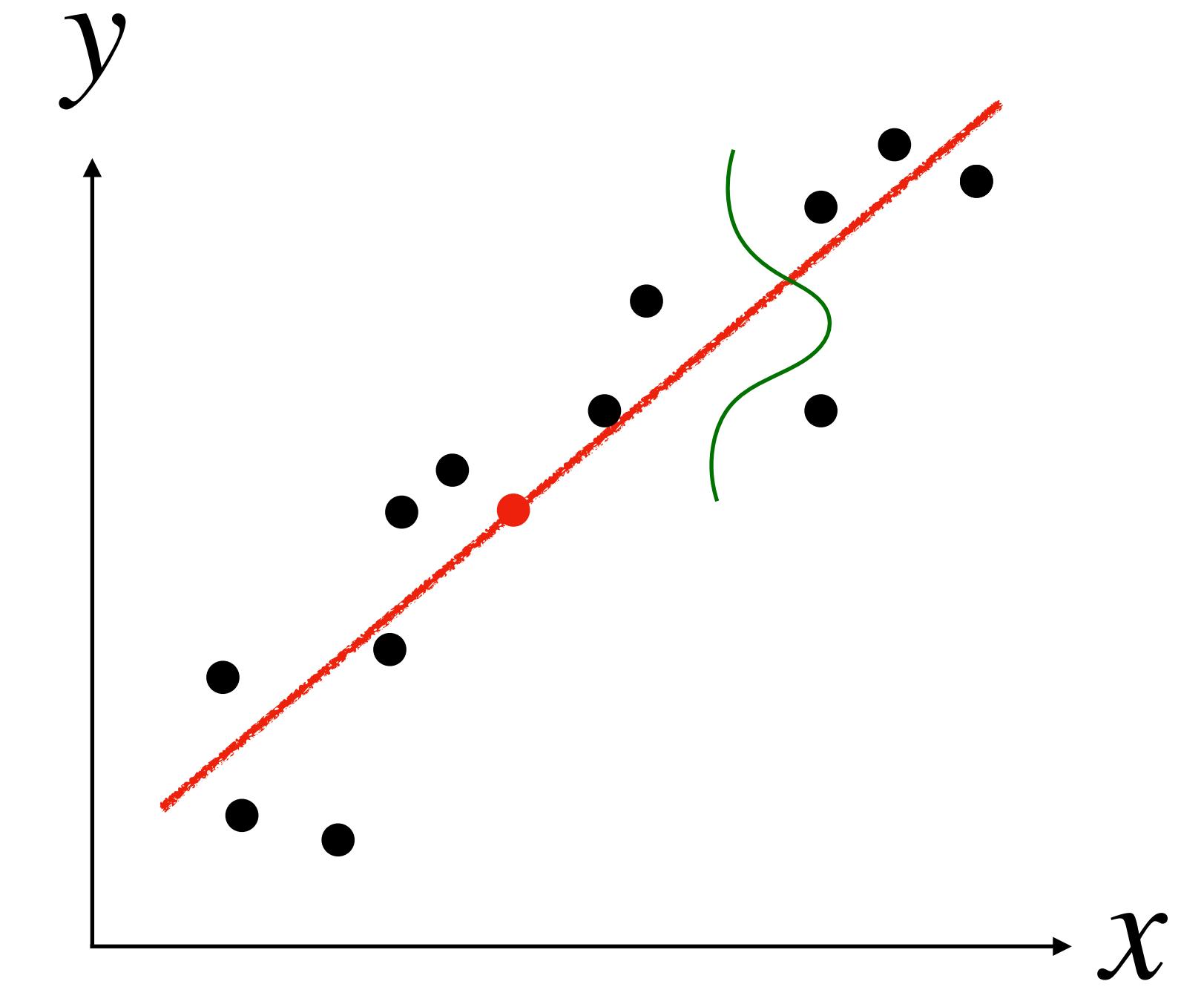
$$L(\theta) = \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta)$$

Independent and Identically Distributed (IID)

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^\top x^{(i)})^2}{2\sigma^2}\right)$$

Log-likelihood

$$\mathcal{L}(\theta) = \log L(\theta)$$
$$= \log \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^\top x^{(i)})^2}{2\sigma^2}\right)$$
$$= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^\top x^{(i)})^2}{2\sigma^2}\right) = n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^\top x^{(i)})^2$$



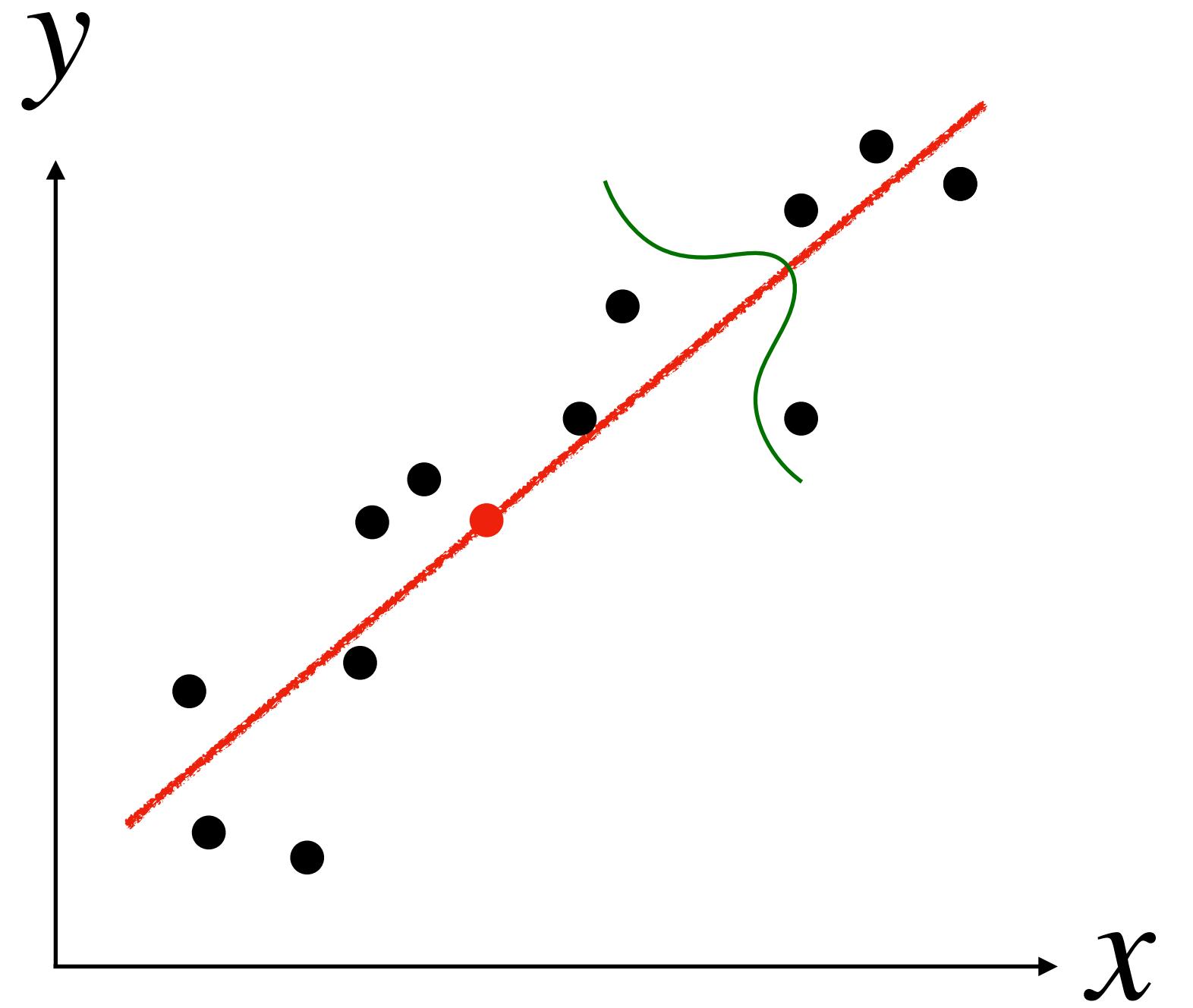
Maximize Log-likelihood

$$\mathcal{L}(\theta) = \log L(\theta)$$

$$= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^\top x^{(i)})^2}{2\sigma^2}\right)$$

$$= n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^\top x^{(i)})^2$$

Maximize $\mathcal{L}(\theta)$ \longrightarrow **Minimize** $\frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^\top x^{(i)})^2$

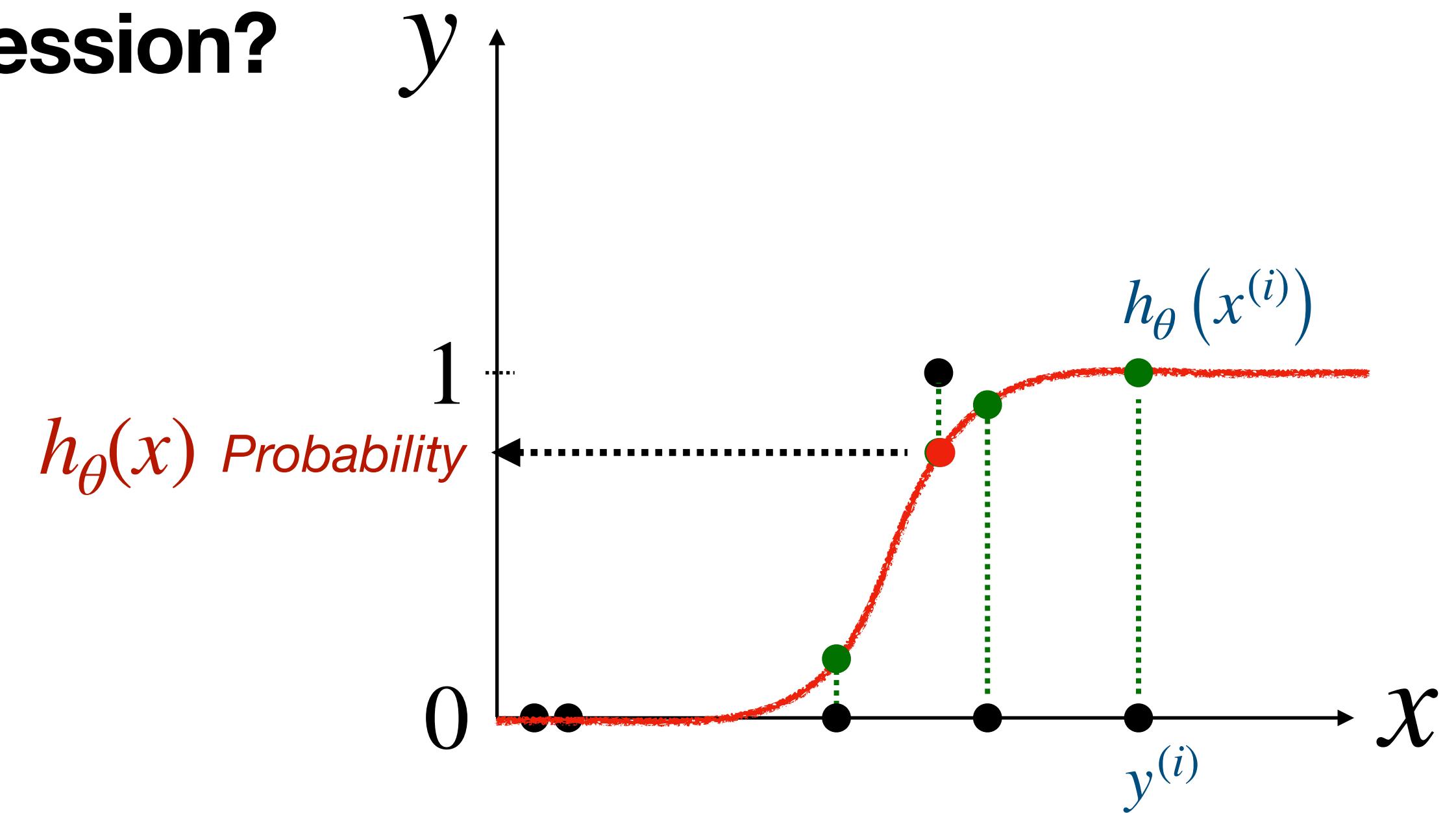


What if the noise is not Gaussian?

Why not Least Squares in Logistic Regression?

$$y = h_{\theta}(x)$$

$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta^T x)}} = \sigma(\theta^T x)$$



Probability of output given input

$$P(y = 1 | x; \theta) = h_{\theta}(x)$$

$$P(y = 0 | x; \theta) = 1 - h_{\theta}(x)$$



$$p(y | x; \theta) = (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y}$$

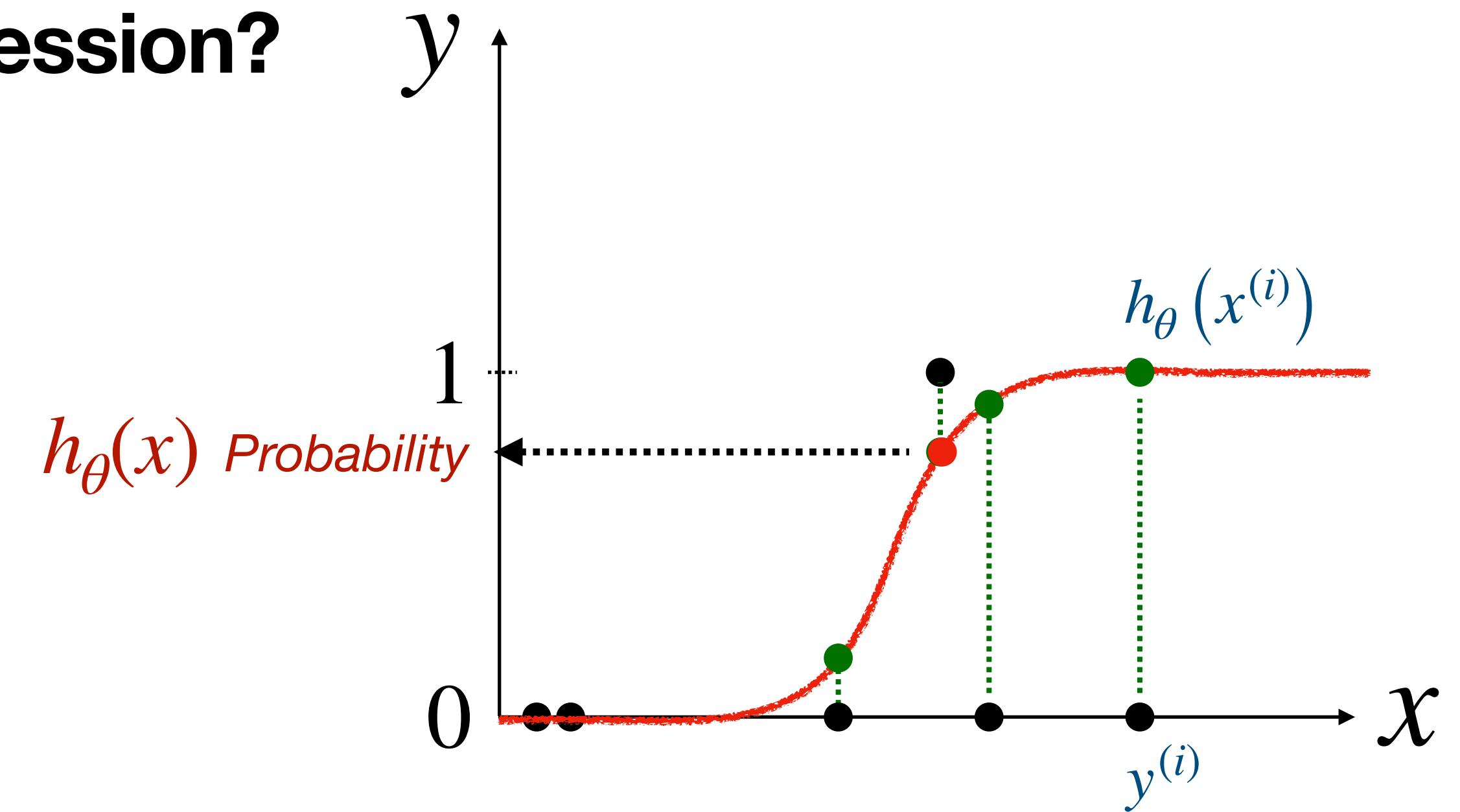
Likelihood!

For Bernoulli Distributed Noise

Why not Least Squares in Logistic Regression?

$$y = h_{\theta}(x)$$

$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta^T x)}} = \sigma(\theta^T x)$$



Probability of output given input

$$P(y = 1 | x; \theta) = \sigma(\theta^T x)$$

$$P(y = 0 | x; \theta) = 1 - \sigma(\theta^T x)$$



$$p(y | x; \theta) = (\sigma(\theta^T x))^y (1 - \sigma(\theta^T x))^{1-y}$$

True label
Likelihood!

For Bernoulli Distributed Noise

Bernoulli Distribution

Properties [edit]

If X is a random variable with a Bernoulli distribution, then:

$$\Pr(X = 1) = p = 1 - \Pr(X = 0) = 1 - q.$$

The probability mass function f of this distribution, over possible outcomes k , is

$$f(k; p) = \begin{cases} p & \text{if } k = 1, \\ q = 1 - p & \text{if } k = 0. \end{cases}$$

[3]

This can also be expressed as

$$f(k; p) = p^k (1 - p)^{1-k} \quad \text{for } k \in \{0, 1\}$$

or as

$$f(k; p) = pk + (1 - p)(1 - k) \quad \text{for } k \in \{0, 1\}.$$

The Bernoulli distribution is a special case of the binomial distribution with $n = 1$.^[4]

Define Log-likelihood

Likelihood

$$p(y|x; \theta) = (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y} \quad \text{for all } (x, y) \text{ pair}$$

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n p(y^{(i)}|x^{(i)}; \theta) \\ &= \prod_{i=1}^n h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}} \end{aligned}$$

$$\log \left(\sum_{i=1}^n y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right)$$

Maximize Log-likelihood

$$\mathcal{L}(\theta) = \sum_{i=1}^n y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)}))$$

Update rule?

while not converged:

$$\theta_j := \theta_j + \alpha \frac{\partial \mathcal{L}(\theta)}{\partial \theta_j}$$

Derive

Gradient Descent

for t = 1...T:

for all parameters j:

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

Maximize Log-likelihood

$$\mathcal{L}(\theta) = \sum_{i=1}^n y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)}))$$

Update rule

while not converged:

$$\theta := \theta + \alpha \nabla_\theta \mathcal{L}(\theta)$$

Derive

Gradient Descent

for t = 1...T:

$$\theta := \theta - \alpha \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$



Same as linear regression



Generalized Linear Models

Gaussian Distribution



Linear Regression

Bernoulli Distribution



Logistic Regression

Update rule

$$\theta := \theta - \alpha \sum_{i=1}^n \left(h_\theta(x^{(i)}) - y^{(i)} \right) x^{(i)}$$

Are there other distributions that lead to the same update rule?

Exponential Family

Family of distributions for which we can derive **the same update rule**

Assumption: $p(y | x; \theta)$ is an exponential family

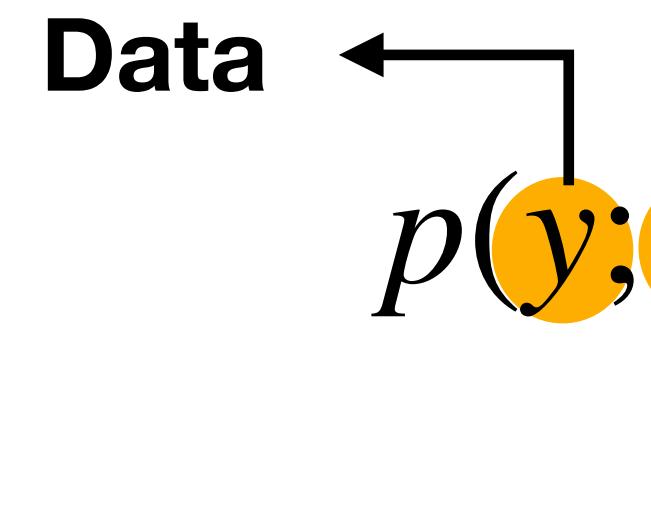
$$p(y; \eta) = b(y) \exp \{ \eta^T y - a(\eta) \}$$

A diagram showing two arrows pointing to the function $p(y; \eta)$. One arrow from the left points to the variable y , labeled "Data". Another arrow from the bottom points to the variable η , labeled "Parameters".

- $b(y)$ is called the base measure (not depend on η)
- $a(\eta)$ is called the log partition function (not depend on y)
- $a(\eta)$, y and $b(y)$ are scalar. η and y have the same dimensions.

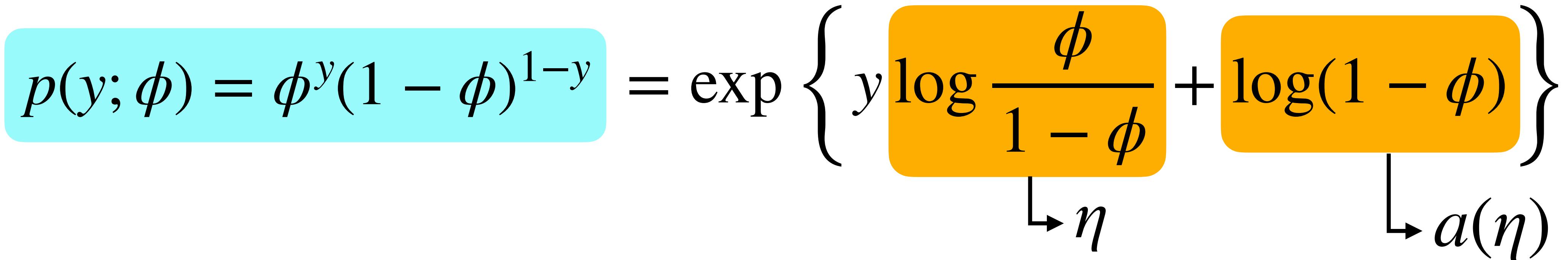
Example 1: Bernoulli Distribution -> Logistic Regression

$$p(y; \eta) = b(y) \exp \left\{ \eta^\top y - a(\eta) \right\}$$

Data  Natural Parameters

Bernoulli Distribution

$$p(y; \phi) = \phi^y (1 - \phi)^{1-y} = \exp \left\{ y \log \frac{\phi}{1 - \phi} + \log(1 - \phi) \right\}$$



Show that term
is only a function of η

Example 2: Gaussian Distribution -> Linear Regression

$$p(y; \eta) = b(y) \exp \left\{ \eta^\top y - a(\eta) \right\}$$

Data ←
→ Natural Parameters

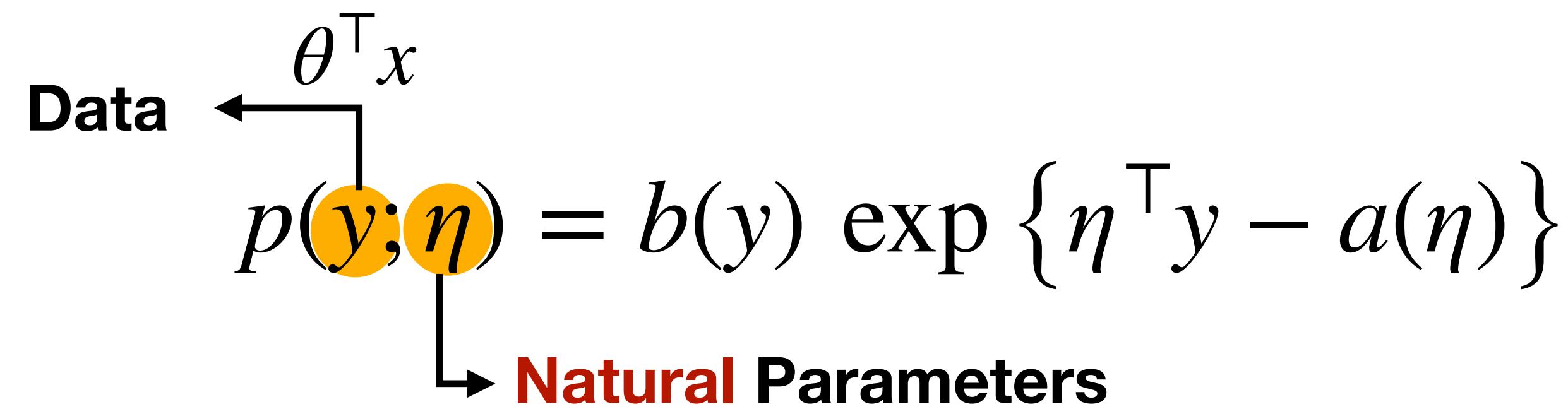
Gaussian Distribution

$$p(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2}(y - \mu)^2 \right\} = \frac{1}{\sqrt{2\pi}} e^{-y^2/2} \left\{ \boxed{\mu y} - \boxed{\frac{1}{2}\mu^2} \right\}$$

$b(y)$ η $a(\eta)$

Why do we care?

$$p(y; \eta) = b(y) \exp \left\{ \eta^\top y - a(\eta) \right\}$$



Inference is Easy:

$$E[y; \eta] = \frac{da(\eta)}{d\eta}$$

$$Var[y; \eta] = \frac{d^2a(\eta)}{d\eta^2}$$

Learning is Easy:

Maximum Likelihood Estimation leads to **convex** problem in η

Generalized Linear Models

Assumption: $p(y | x; \theta)$ is an exponential family

Data Type → Probability Distribution

Binary → Bernoulli → **Logistic Regression**

Real → Gaussian → **Linear Regression**

Counts → Poisson

Positive Real → Gamma, Exponential

Distributions → Dirichlet

Generalized Linear Models

Assumption: $p(y | x; \theta)$ is an exponential family

The natural parameter is linear in the inputs

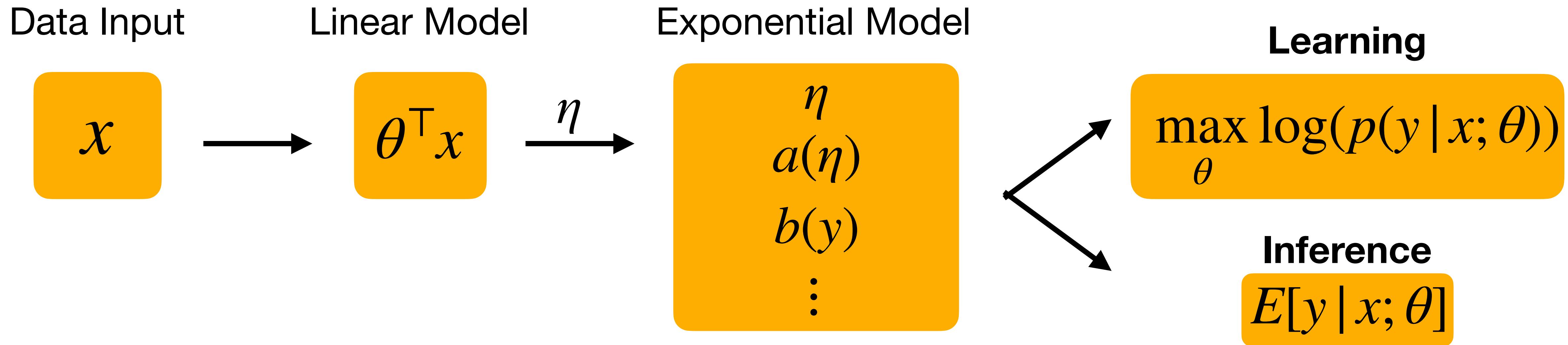
$$\eta = \theta^T x$$

Predictor is a natural consequence

$$h_\theta(x) = E[y | x; \theta]$$

Generalized Linear Models

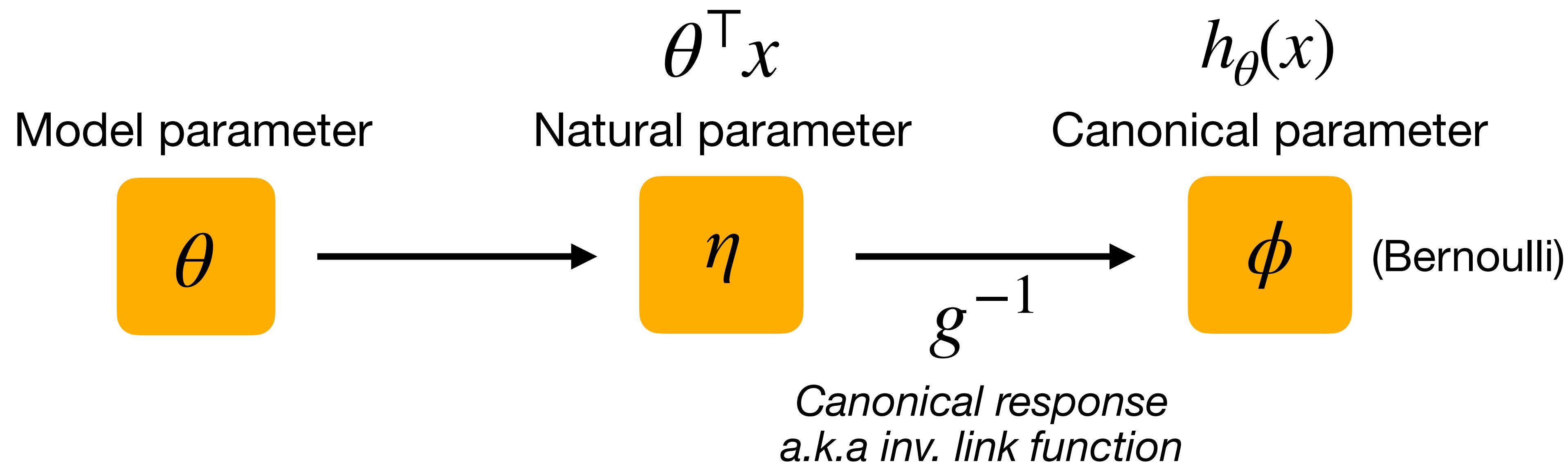
Assumption: $p(y | x; \theta)$ is an exponential family



Update Rule:

$$\theta := \theta - \alpha \sum_{i=1}^n \left(h_{\theta} (x^{(i)}) - y^{(i)} \right) x^{(i)}$$

Terminology



Logistic Regression:

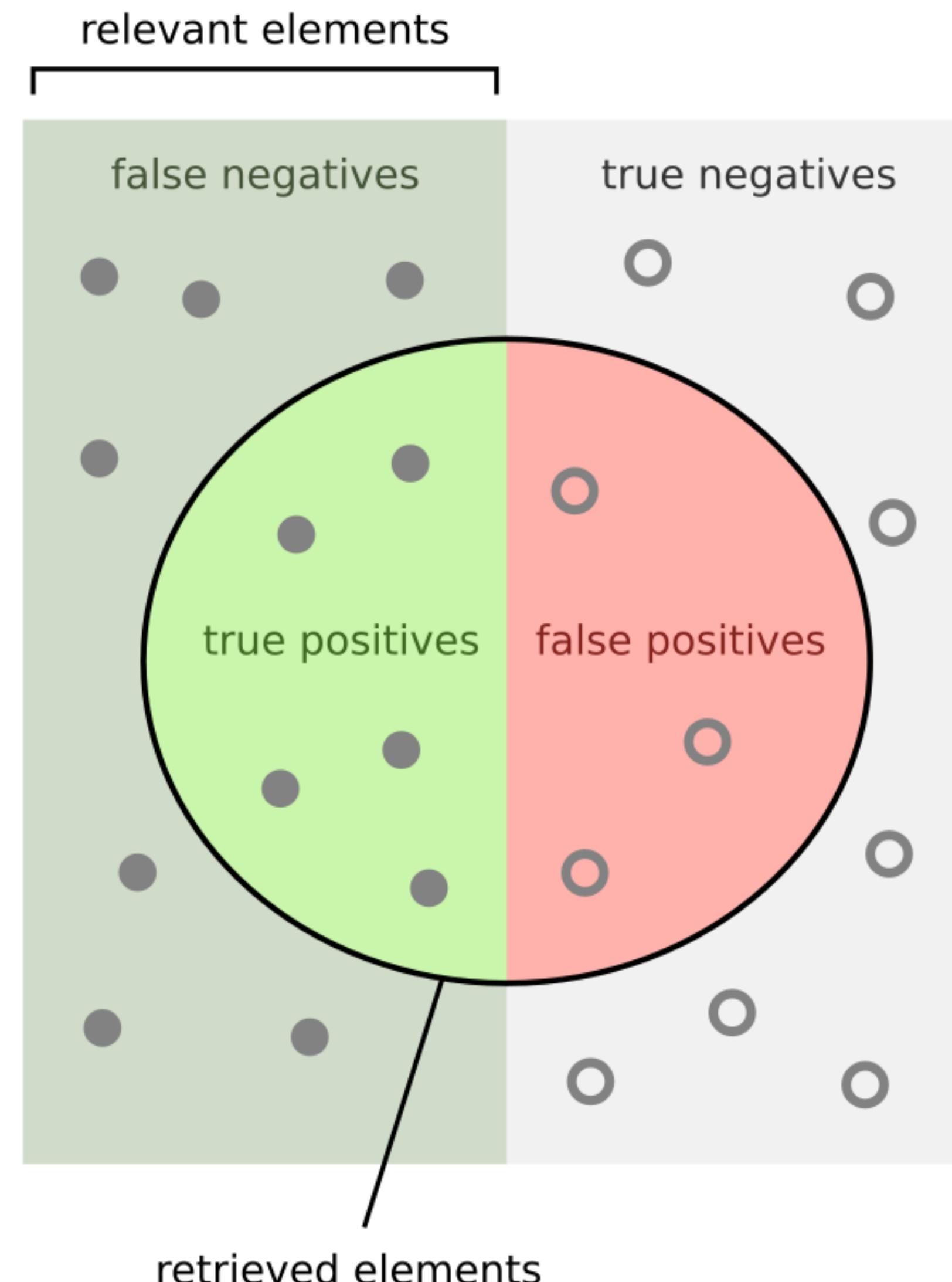
$$h_\theta(x) = E[y | x; \theta]$$

$$\phi = \frac{1}{1 + e^{-\eta}} = \frac{1}{1 + e^{-\theta^\top x}}$$

Evaluation Metrics

		Predicted condition		Sources: [4][5][6][7][8][9][10][11] view · talk · edit		
		Total population $= P + N$	Predicted positive	Predicted negative	Informedness, bookmaker informedness (BM) $= TPR + TNR - 1$	Prevalence threshold $\frac{(PT)}{\sqrt{TPR \times FPR} - FPR} = \frac{TPR - FPR}{TPR - FPR}$
Actual condition	Positive (P) [a]	True positive (TP), hit ^[b]	False negative (FN), miss, underestimation	True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power $= \frac{TP}{P} = 1 - FNR$	False negative rate (FNR), miss rate type II error [c] $= \frac{FN}{P} = 1 - TPR$	
	Negative (N) ^[d]	False positive (FP), false alarm, overestimation	True negative (TN), correct rejection ^[e]	False positive rate (FPR), probability of false alarm, fall-out type I error [f] $= \frac{FP}{N} = 1 - TNR$	True negative rate (TNR), specificity (SPC), selectivity $= \frac{TN}{N} = 1 - FPR$	
Prevalence $= \frac{P}{P + N}$	Positive predictive value (PPV), precision $= \frac{TP}{TP + FP} = 1 - FDR$	Negative predictive value (NPV) $= \frac{TN}{TN + FN} = 1 - FOR$	Positive likelihood ratio (LR+) $= \frac{TPR}{FPR}$	Negative likelihood ratio (LR-) $= \frac{FNR}{TNR}$		
Accuracy (ACC) $= \frac{TP + TN}{P + N}$	False discovery rate (FDR) $= \frac{FP}{TP + FP} = 1 - PPV$	False omission rate (FOR) $= \frac{FN}{TN + FN} = 1 - NPV$	Markedness (MK), deltaP (Δp) $= PPV + NPV - 1$	Diagnostic odds ratio (DOR) $= \frac{LR+}{LR-}$		
Balanced accuracy (BA) $= \frac{TPR + TNR}{2}$	F_1 score $= \frac{2 \cdot PPV \times TPR}{PPV + TPR}$ $= \frac{2 \cdot TP}{2 \cdot TP + FP + FN}$	Fowlkes–Mallows index (FM) $= \sqrt{PPV \times TPR}$	phi or Matthews correlation coefficient (MCC) $= \sqrt{TPR \times TNR \times PPV \times NPV} - \sqrt{FNR \times FPR \times FOR \times DFR}$	Threat score (TS), critical success index (CSI), Jaccard index $= \frac{TP}{TP + FN + FP}$		

Precision and Recall



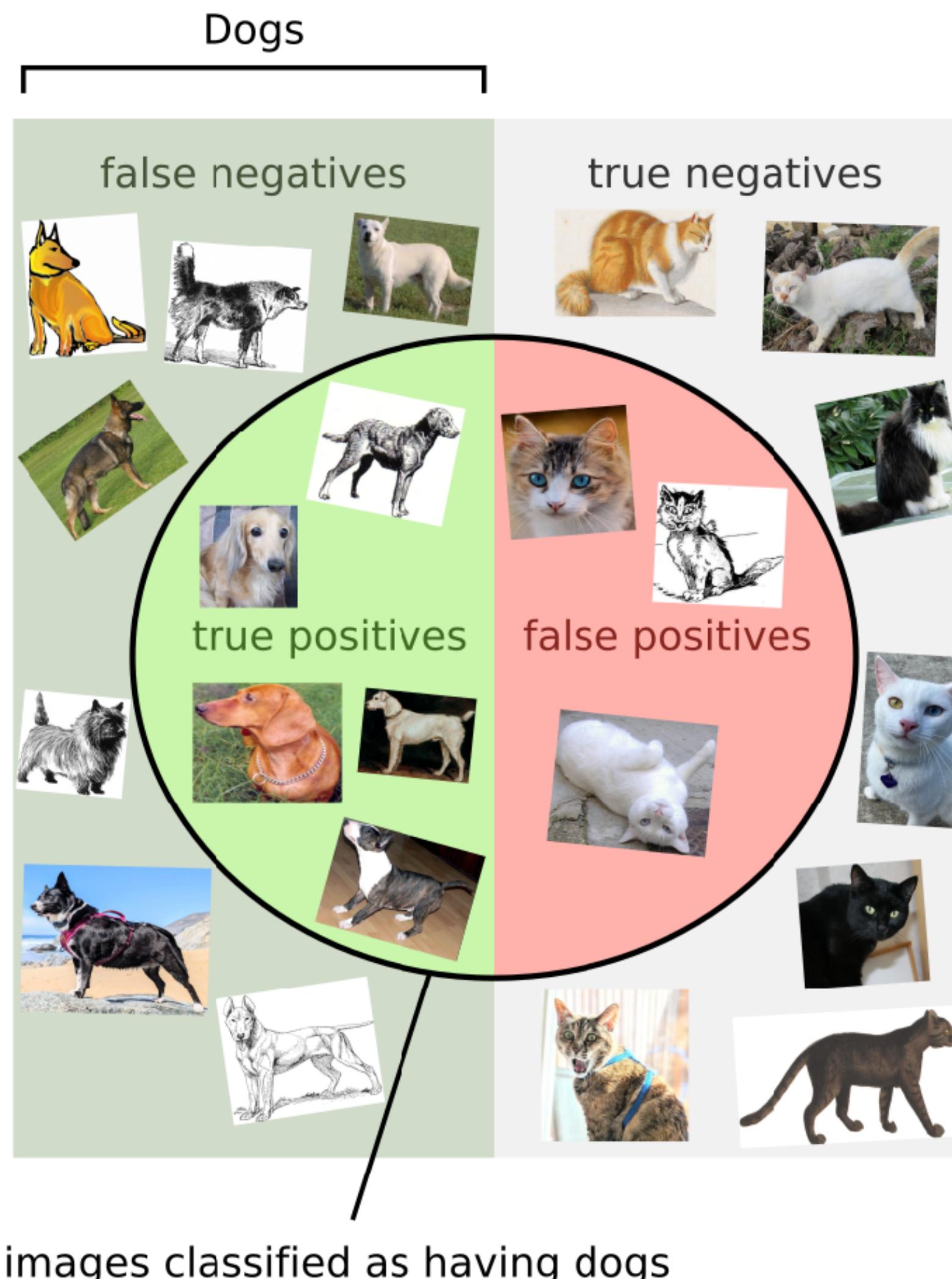
How many retrieved items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant items are retrieved?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

Precision and Recall



$$\text{Precision} = \frac{5 \text{ true pos.}}{8 \text{ total pos.}}$$

$$\text{Recall} = \frac{5 \text{ true pos.}}{12 \text{ total dogs}}$$

$$\text{Prevalence} = \frac{12 \text{ total dogs}}{22 \text{ total images}}$$

$$\text{Accuracy} = \frac{5 \text{ true pos.} + 7 \text{ true neg.}}{22 \text{ total images}}$$

$$F_1 = \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}} = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$$

```
sklearn.metrics.f1_score(y_true, y_pred, *, labels=None, pos_label=1,
average='binary', sample_weight=None, zero_division='warn') # [source]
```

Compute the F1 score, also known as balanced F-score or F-measure.

The F1 score can be interpreted as a harmonic mean of the precision and recall, where an F1 score reaches its best value at 1 and worst score at 0. The relative contribution of precision and recall to the F1 score are equal. The formula for the F1 score is:

Accuracy and Precision

