

EXPONENTIAL FAMILIES & GLMs

$$y|x; \theta \sim N(\boxed{\theta^T x}, \sigma^2) \quad \text{Regression } y \in \mathbb{R}$$

$$y|x; \theta \sim \text{Bern}\left(\frac{1}{1 + e^{-\boxed{\theta^T x}}}\right) \quad \text{Classification } y \in \{0, 1\}$$

pdf/
pmf

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

most
cases
↓

$b(y)$ - Base Measure

\boxed{y}

$T(y)$ - Sufficient Statistic

$a(\eta)$ - log partition function
 $\eta \rightarrow$ natural parameter
probability distribution

$$p(y; \eta) \propto b(y) e^{\eta^T y}$$

$$p(y; \eta) = \frac{b(y) e^{\eta^T y}}{\int b(y) e^{\eta^T y} dy} = \frac{b(y) e^{\eta^T y}}{E[e^{\eta^T y}]}$$

\uparrow M.G.F.

$$= b(y) e^{\eta^T y - \log A(\eta)}$$

$$\boxed{a(\eta) = \log A(\eta)}$$

Bernoulli

$$\begin{aligned} p(y; \phi) &= \phi^y (1-\phi)^{1-y} \\ &= \exp[\log(\phi^y (1-\phi)^{1-y})] \\ &= \exp[y \log \phi + (1-y) \log(1-\phi)] \\ &= \exp[y \log \phi - y \log(1-\phi) + \log(1-\phi)] \\ &= \exp\left[y \log\left(\frac{\phi}{1-\phi}\right) + \log(1-\phi)\right] \end{aligned}$$

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

$$\boxed{\eta = \log\left(\frac{\phi}{1-\phi}\right)}$$

$$\boxed{T(y) = y}$$

$$\phi = \frac{1}{1+e^{-\eta}} \rightarrow \text{Logistic function}$$

$$a(\eta) = -\log(1-\phi) \\ = \log(1+e^\eta)$$

$$b(y) = 1$$

Gaussian

$$P(y; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2} \frac{(y-\mu)^2}{\sigma^2}\right\}$$

Assume $\sigma^2 = 1$

$$P(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}(y-\mu)^2\right\} \\ = \left[\frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}y^2\right\} \right] \exp\left\{\mu y - \frac{1}{2}\mu^2\right\}$$

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

$$\eta = \mu \quad a(\eta) = +\frac{1}{2}\mu^2$$

$$T(y) = y \quad = \frac{1}{2}\eta^2$$

$$b(y) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}y^2\right\}$$

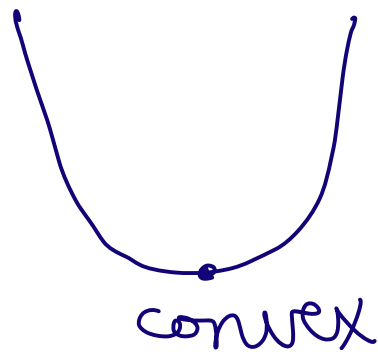
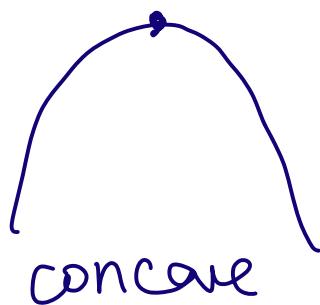
↳ Standard
Gaussian
 $\mu = 0$
 $\sigma = 1$

<1> $\log P(y; \eta)$ is concave in η
 $MLE \rightarrow$ concave in η
 $NLL \rightarrow$ convex in η .

$$<2> E[y; \eta] = \frac{\partial}{\partial \eta} a(\eta)$$

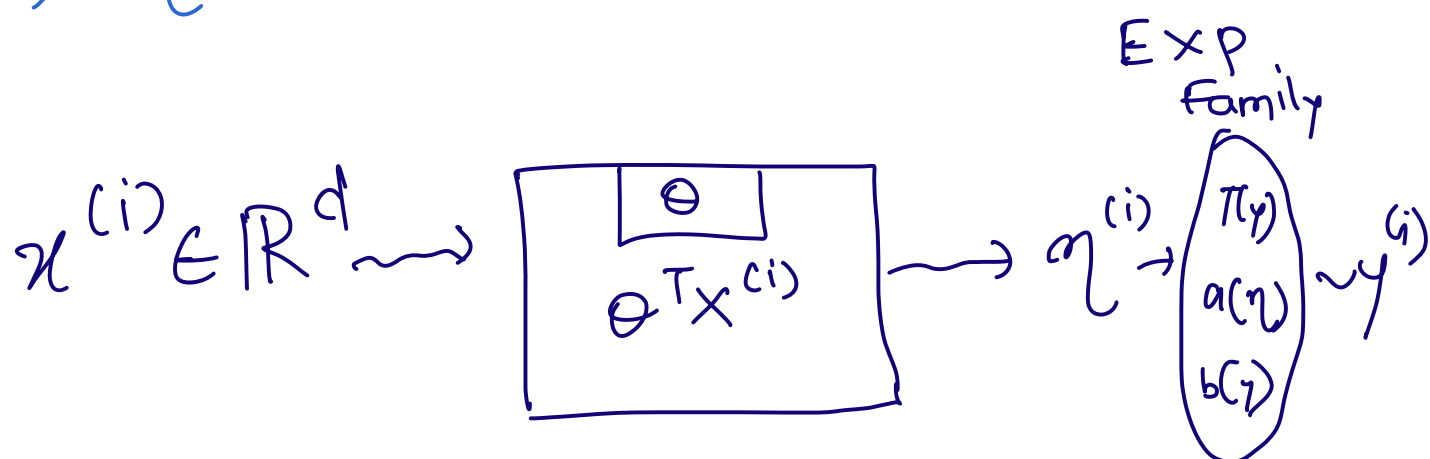
$$<3> \text{var}[y; \eta] = \frac{\partial^2 a(\eta)}{\partial^2 \eta}$$

These are especially nice because when given distribution, to find E , you have to integrate but here you have differentiation.



G.L.M.

- <1> $y|x; \theta \sim \text{ExpFamily}(\eta)$
- <2> $h_\theta(x) = E[y|x; \theta]$
- <3> $\eta = \theta^T x$



We first fix a distribution which
fixes $\pi(y), a(\eta), b(\eta)$ which models
our output

$\theta \rightarrow \text{global}$
 $x^{(i)} \rightarrow \text{local}$
 $\eta = \theta^T x = \text{local}$

LINEAR REGRESSION

Exp Family \rightarrow Gaussian

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

$$= \mu$$

$$= \eta$$

$$= \theta^T x$$

LOGISTIC

Exp family is Bernoulli

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

$$= \phi$$

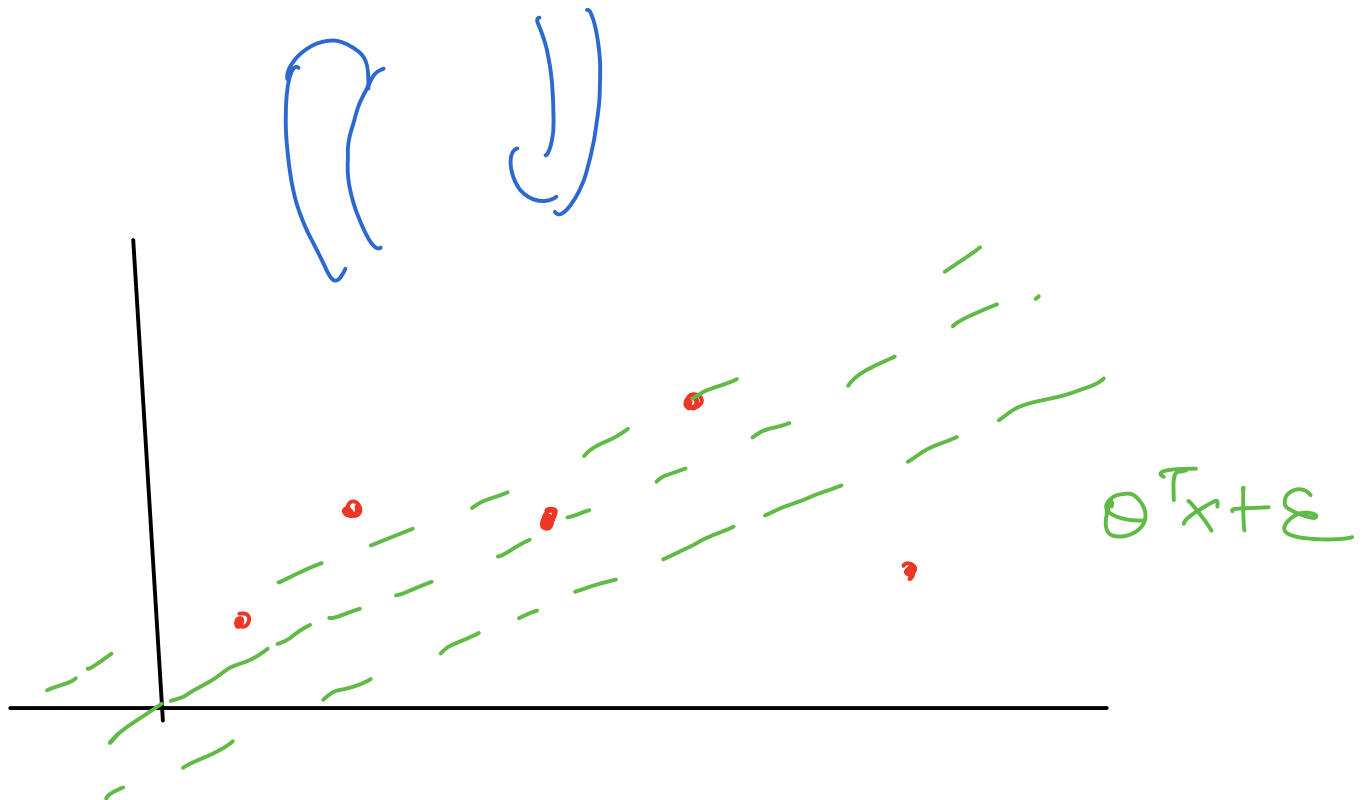
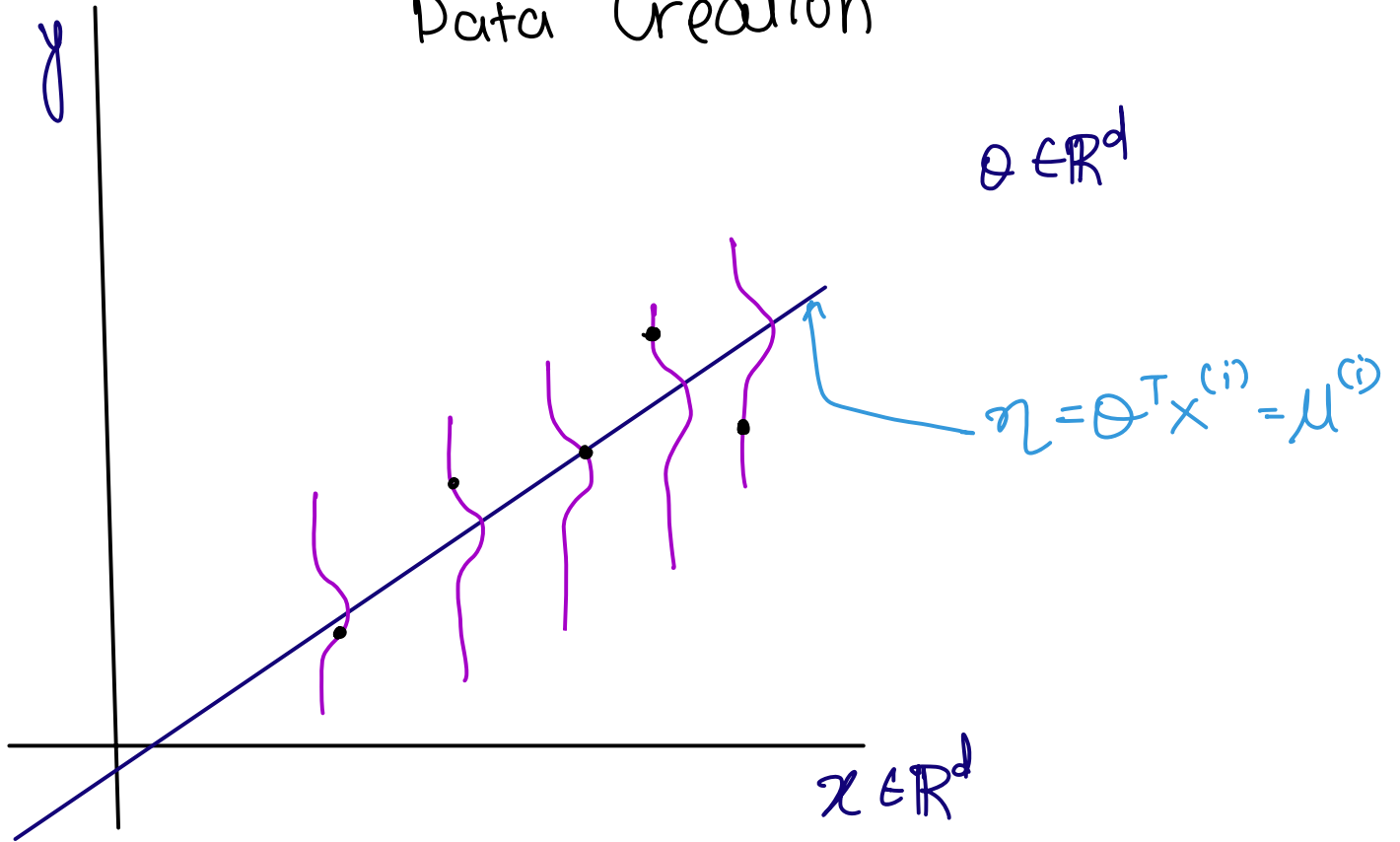
$$= \frac{1}{1+e^{-\eta}}$$

$$= \frac{1}{1+e^{-\theta^T x}}$$

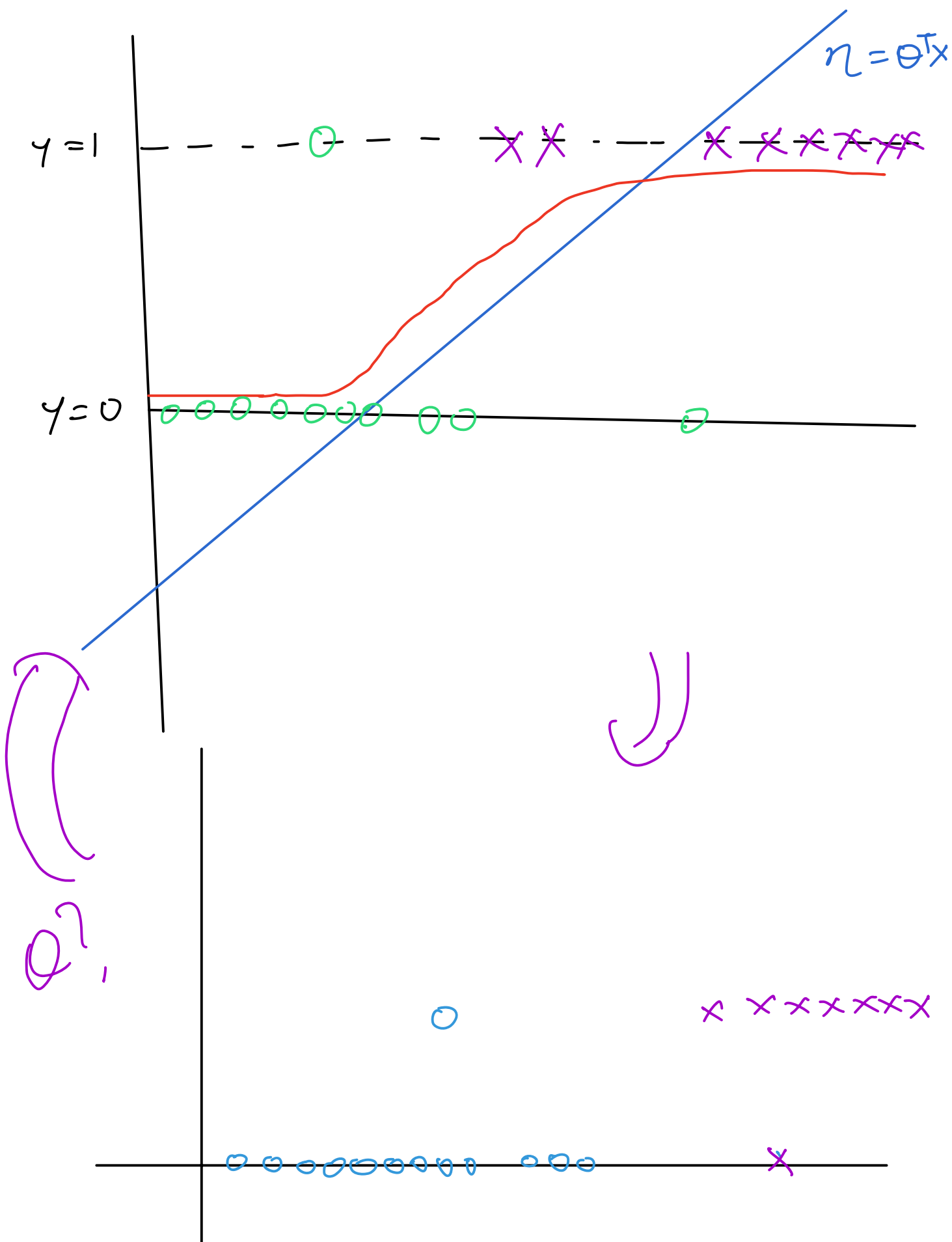
$$= g(\theta^T x)$$

$$g(z) = \frac{1}{1+e^{-z}}$$

Data Creation



We assume first that data generated according to Gaussian. Then when we have data, we actually don't know θ , we assume there is true θ & work to find it $\hat{\theta}$ that is close to true θ .



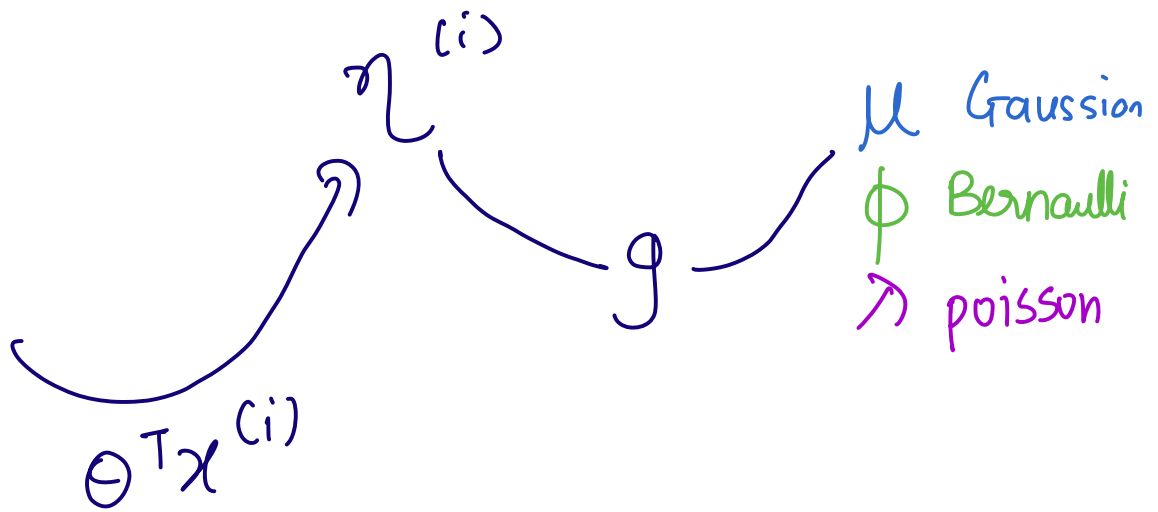
MODEL

NATURAL

MEAN

$$\theta \in \mathbb{R}^d$$

$$x^{(i)} \in \mathbb{R}^d$$



Canonical Response function

Gaussian $g(\eta) = \eta$

Bern $\phi = g(\eta) = \frac{1}{1+e^{-\eta}}$

$g^{-1} \rightarrow$ Canonical link function

$y \rightarrow$ response function

$$\text{Prediction: } h_{\theta}(x) = E[y|x; \theta] = g(\theta^T x)$$

G.L.M.

Data Type	Exp Fam Dist ⁿ	Name
\mathbb{R}	Gaussian Laplace	Regression
$\{0, 1\}$	Bernoulli	Classification
$\{1, \dots, K\}$	Categorical	Multiclass Classification
\mathbb{N}_+	Poisson	Count Reg/ Poisson Reg
$\mathbb{R}_+(\text{time})$	Exponential, Gamma	Survival Analysis

<1> Make choice of dist according to data type

<2> Express in exp. form: $-a(\eta), b(\eta), T(\eta)$
 $\mu, \phi = g(\eta)$

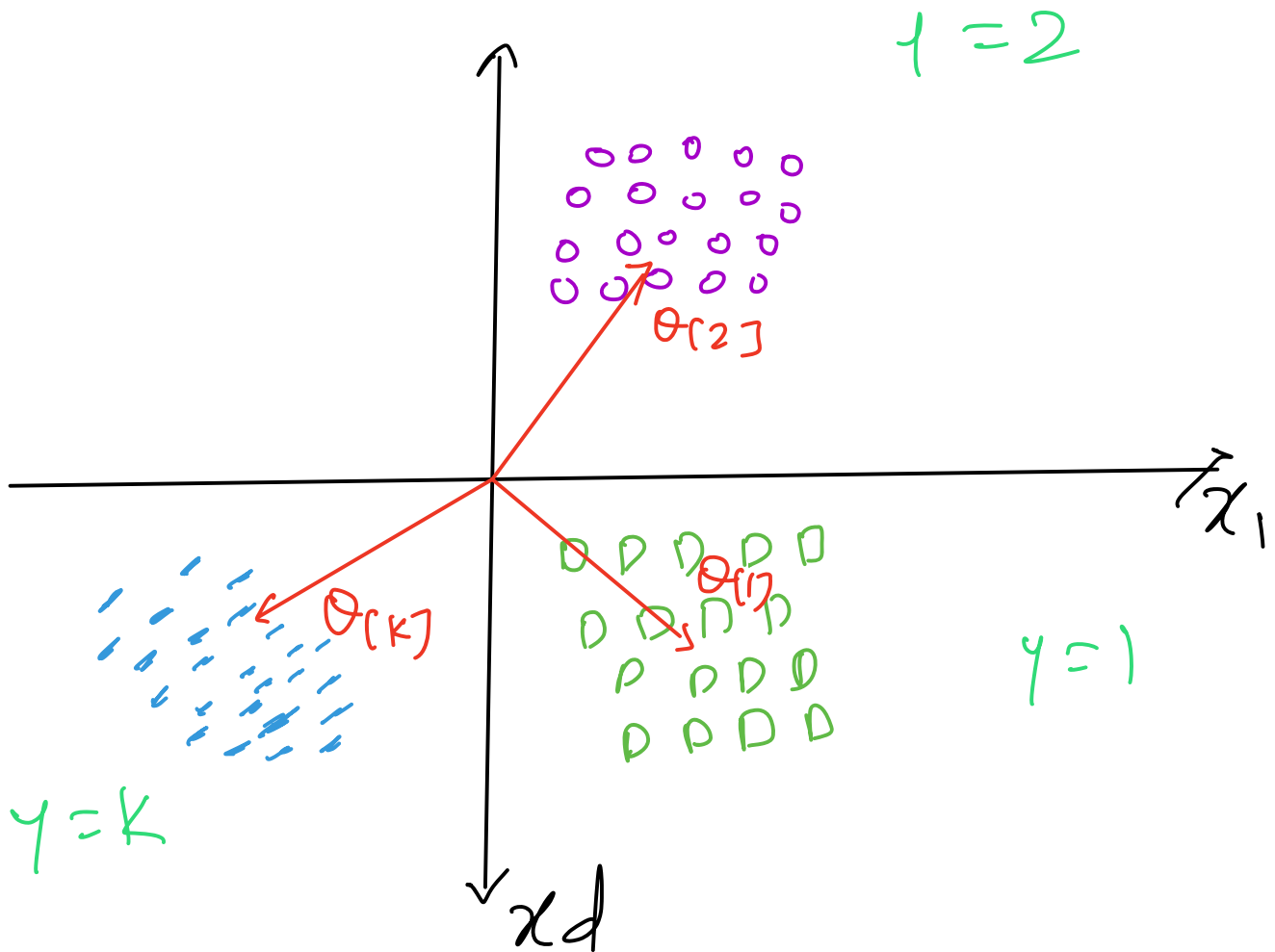
<3> Hypothesis $h_{\theta}(x) = g(\theta^T x)$

<4> $\theta := \theta + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) \cdot x^{(i)}$

<5> Prediction: $\hat{y} = h_{\theta}(\vec{x}) = g(\theta^T \vec{x})$

SOFTMAX REGRESSION

Multiclass Classification



One θ per class

$$\left. \begin{array}{l} \theta_{[1]}^T x - R \\ \theta_{[2]}^T x - R \\ \vdots \\ \theta_{[k]}^T x - R \end{array} \right\} \arg \max_{\theta}$$

$$\theta = \begin{bmatrix} -\theta_{[1]}^T & - \\ \vdots & \\ -\theta_{[k]}^T & - \end{bmatrix} \begin{matrix} \uparrow \\ k \\ \downarrow \end{matrix} \begin{bmatrix} x \\ \vdots \end{bmatrix} \begin{matrix} \uparrow \\ d \\ \downarrow \end{matrix} \begin{matrix} \theta_{[1]}^T x \\ \theta_{[2]}^T x \\ \vdots \\ \theta_{[k]}^T x \end{matrix}$$

Then $\exp \theta^T x$ & normalise it

$$\frac{\exp(x^T \theta_{[i]})}{\sum_j \exp(x^T \theta_{[j]})}$$

Maximum component of \vec{x}

$$\rightarrow \vec{x}^T \text{softmax}(\vec{z}) \approx \max(\vec{z})$$