

Probabilistic Methods-II

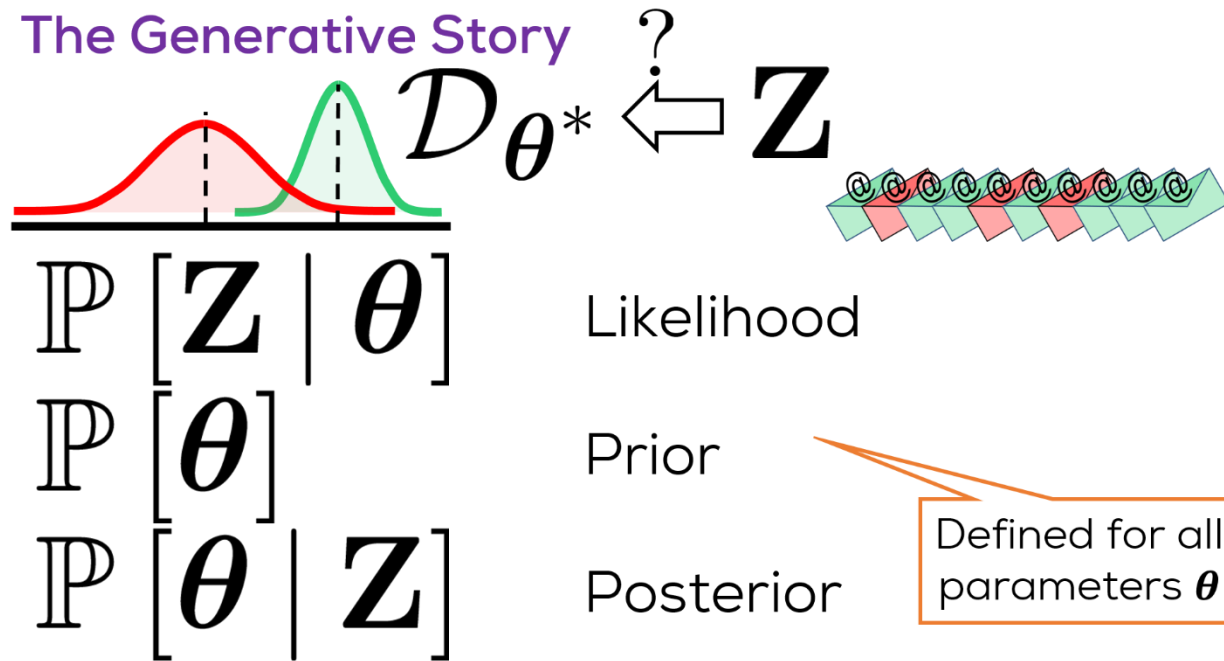
CS771: Introduction to Machine Learning

Purushottam Kar



Recap

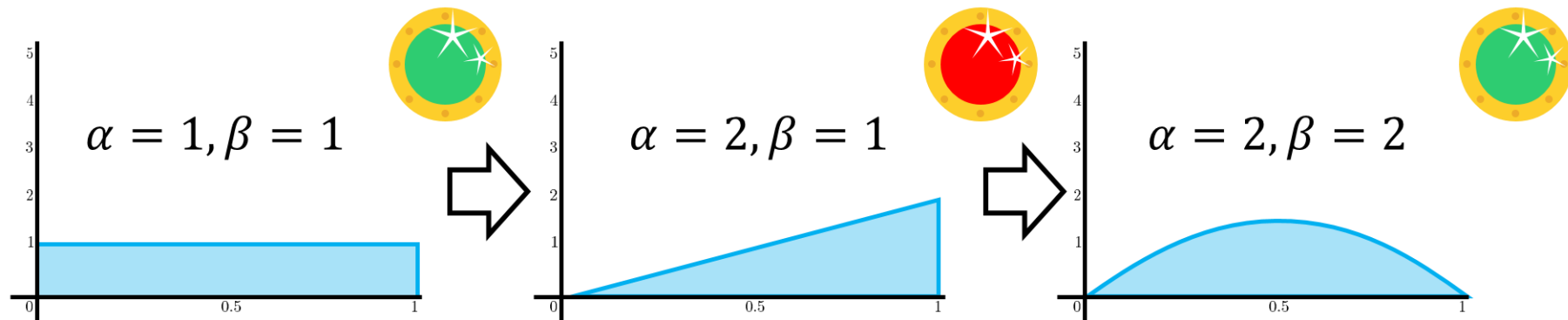
The Generative Story



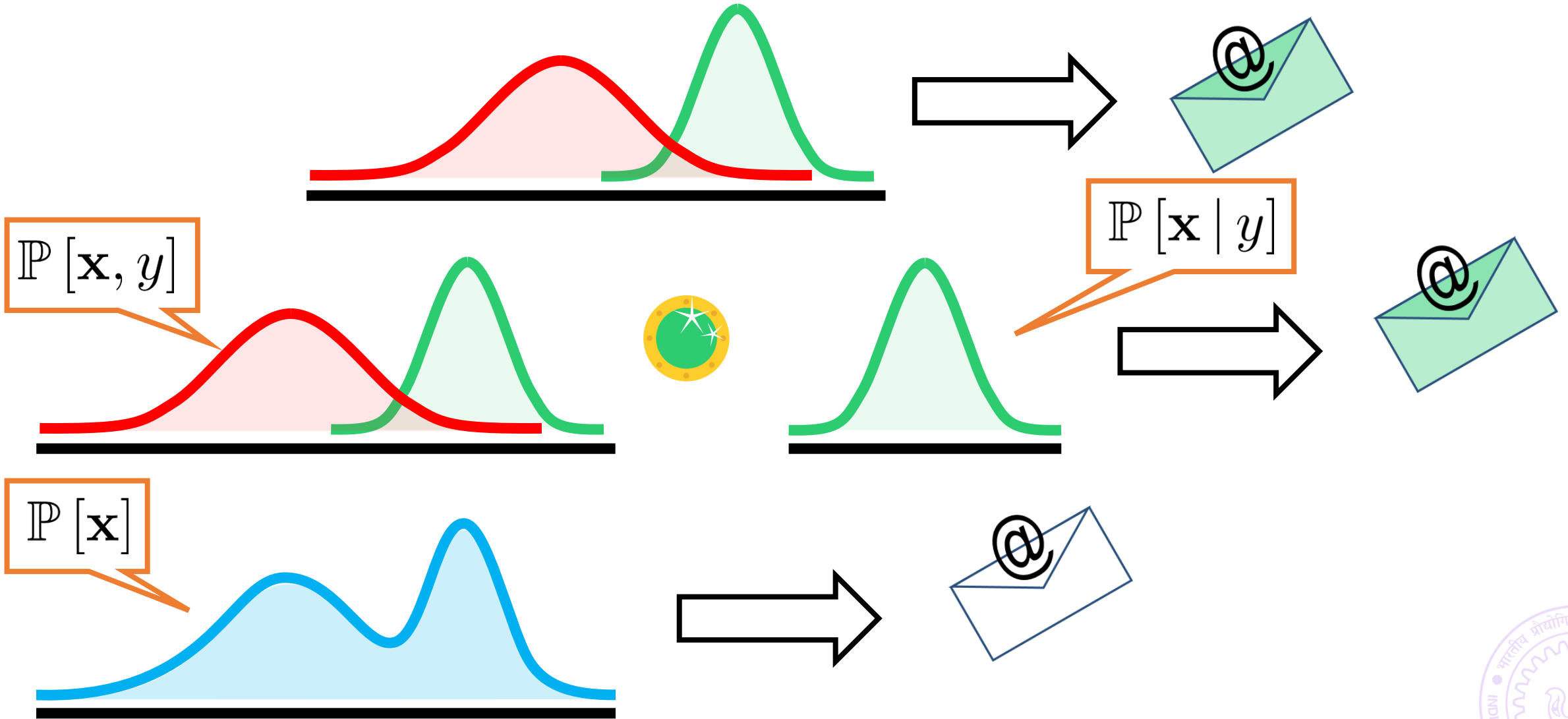
Prediction with learnt models

Diagram illustrating prediction with learnt models. A coin flip is shown with a question mark. The data $\mathbf{Z} = \mathbf{y}$ is represented by a sequence of colored circles (green and red). The probability $\mathbb{P}[p]$ is shown as a horizontal bar. The likelihood is given by $\mathbb{P}[y^i | p] = p^{y^i} (1 - p)^{1 - y^i}$. The posterior probability is given by $\mathbb{P}[\text{coin} = H | \mathbf{y}] = \int_p \mathbb{P}[\text{coin} = H | p] \mathbb{P}[p | \mathbf{y}] dp$. The maximum likelihood estimate (MLE) is $\hat{p}_{\text{MLE}} = \frac{n_H}{n}$. The maximum a posteriori (MAP) estimate is $\hat{p}_{\text{MAP}} = \frac{n_H + \alpha - 1}{n + \alpha + \beta - 2}$.

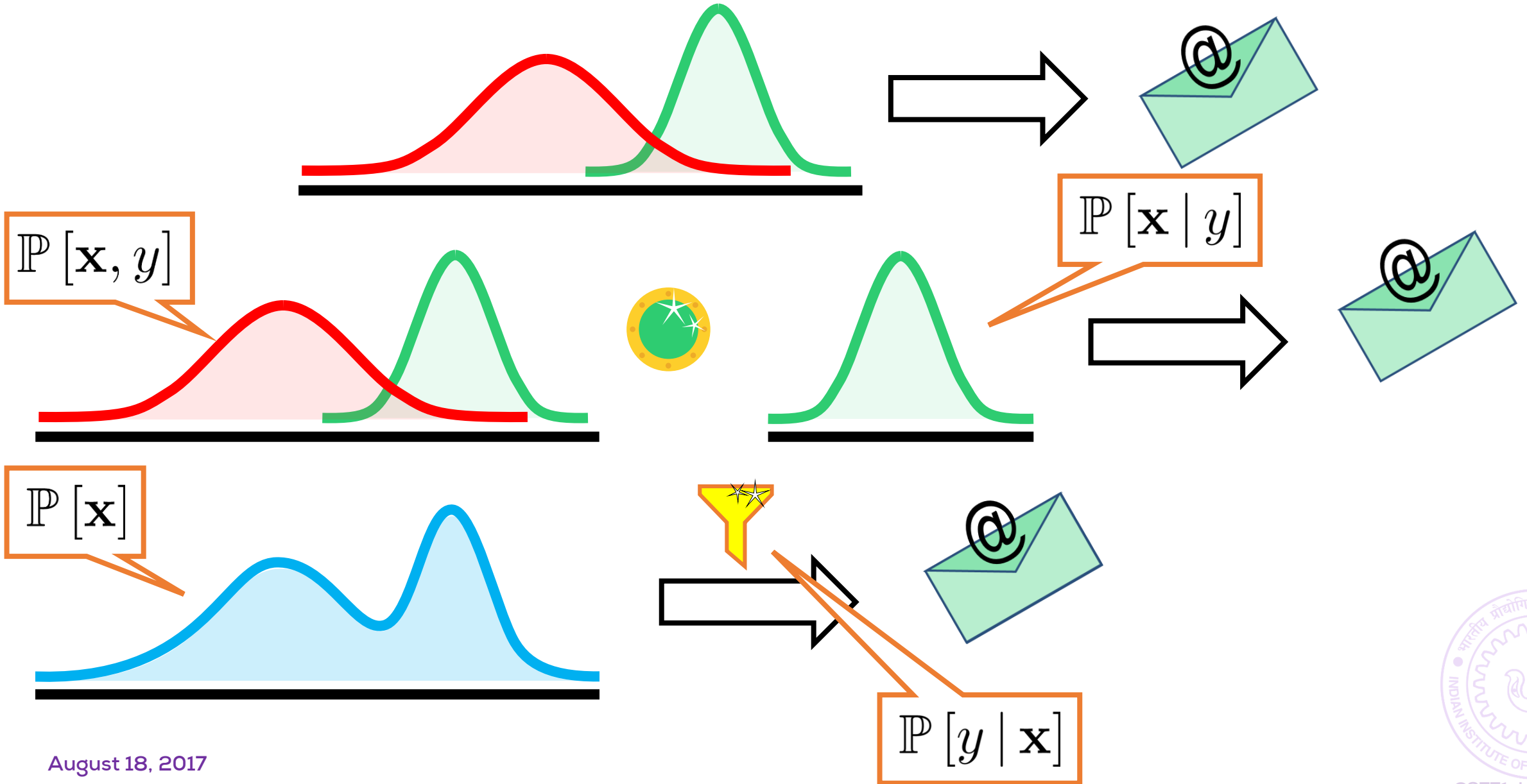
Online MAP!



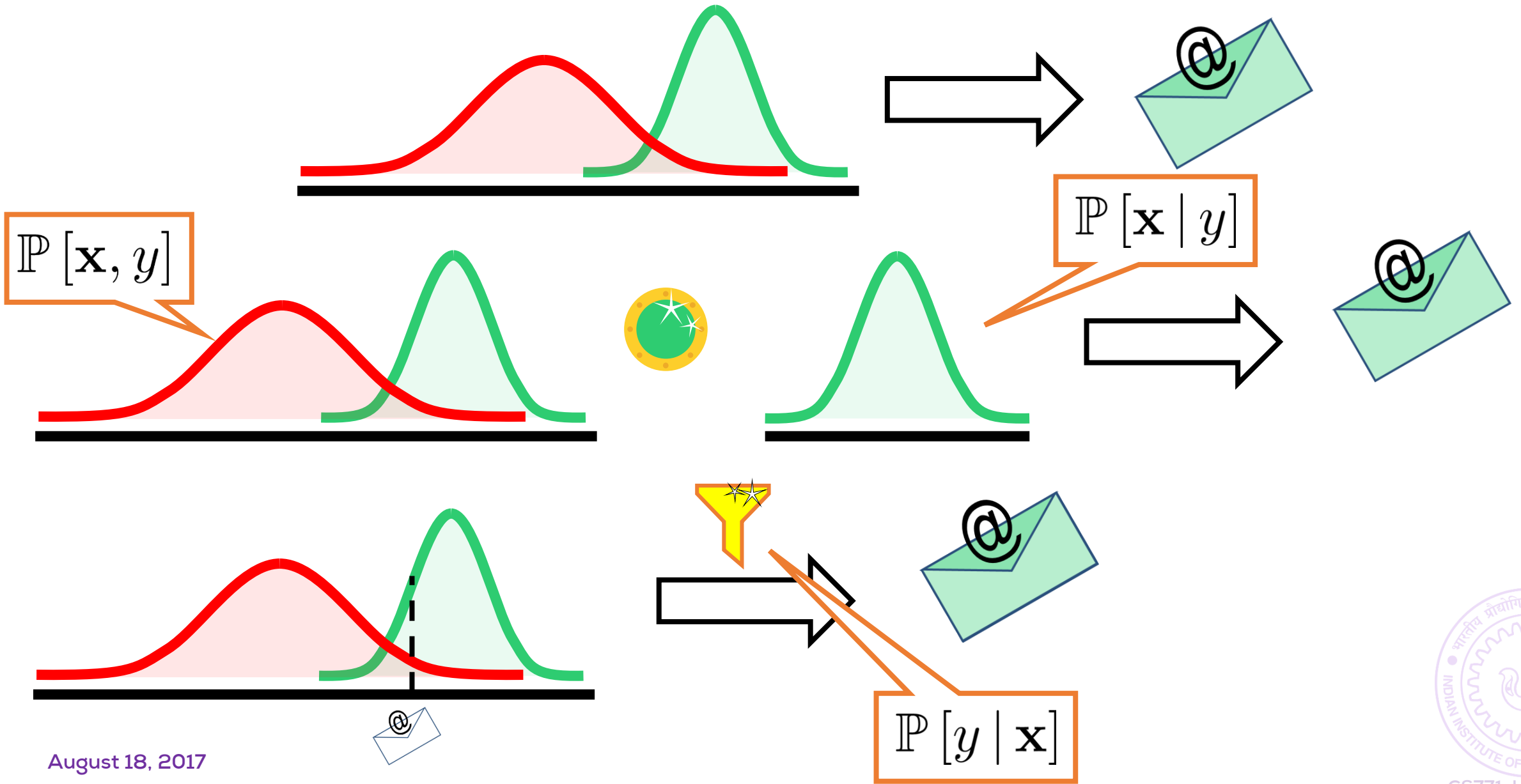
Generative Models for Labeled Data



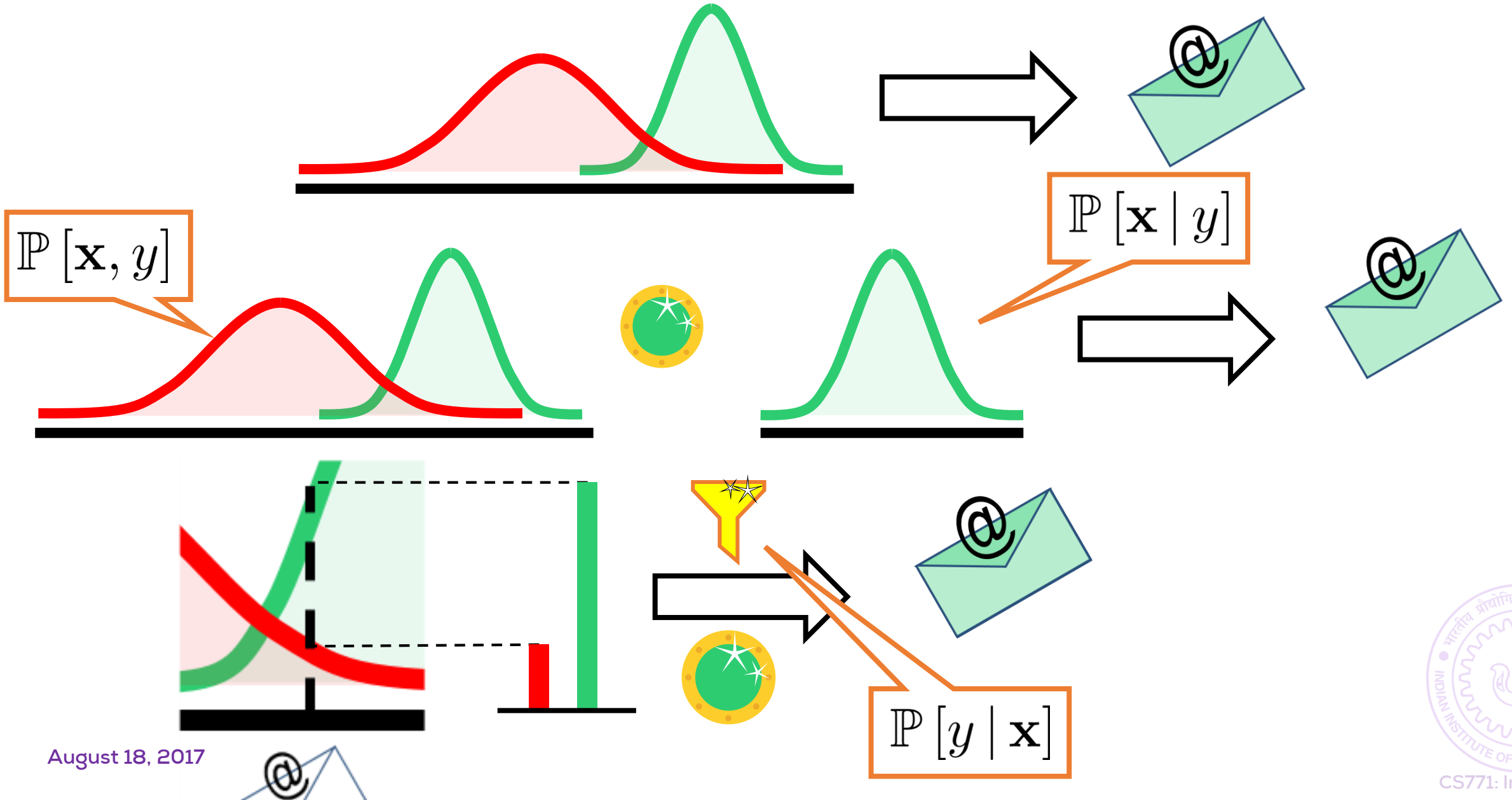
Generative Models for Labeled Data



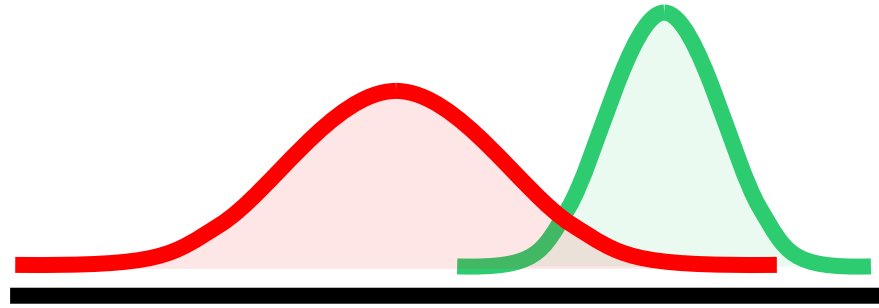
Generative Models for Labeled Data



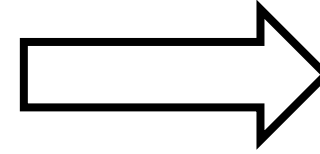
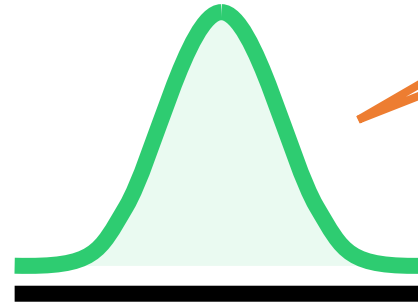
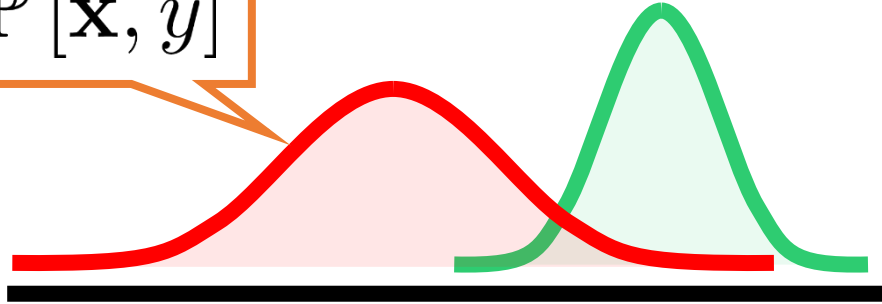
Generative Models for Labeled Data



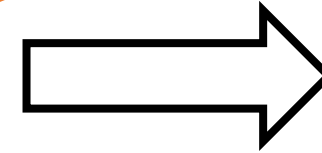
Generative Models for Labeled Data



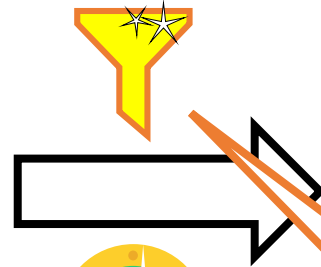
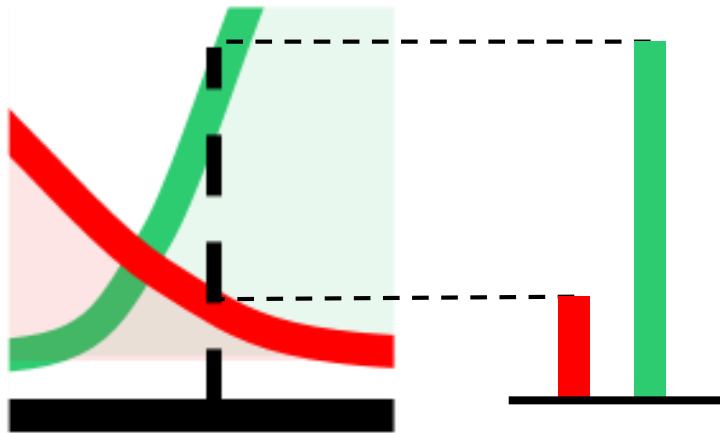
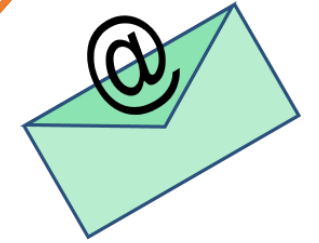
$\mathbb{P}[\mathbf{x}, y]$



$\mathbb{P}[\mathbf{x} | y]$



Used in
generative
learning



Used in
discriminative
learning

$\mathbb{P}[y | \mathbf{x}]$

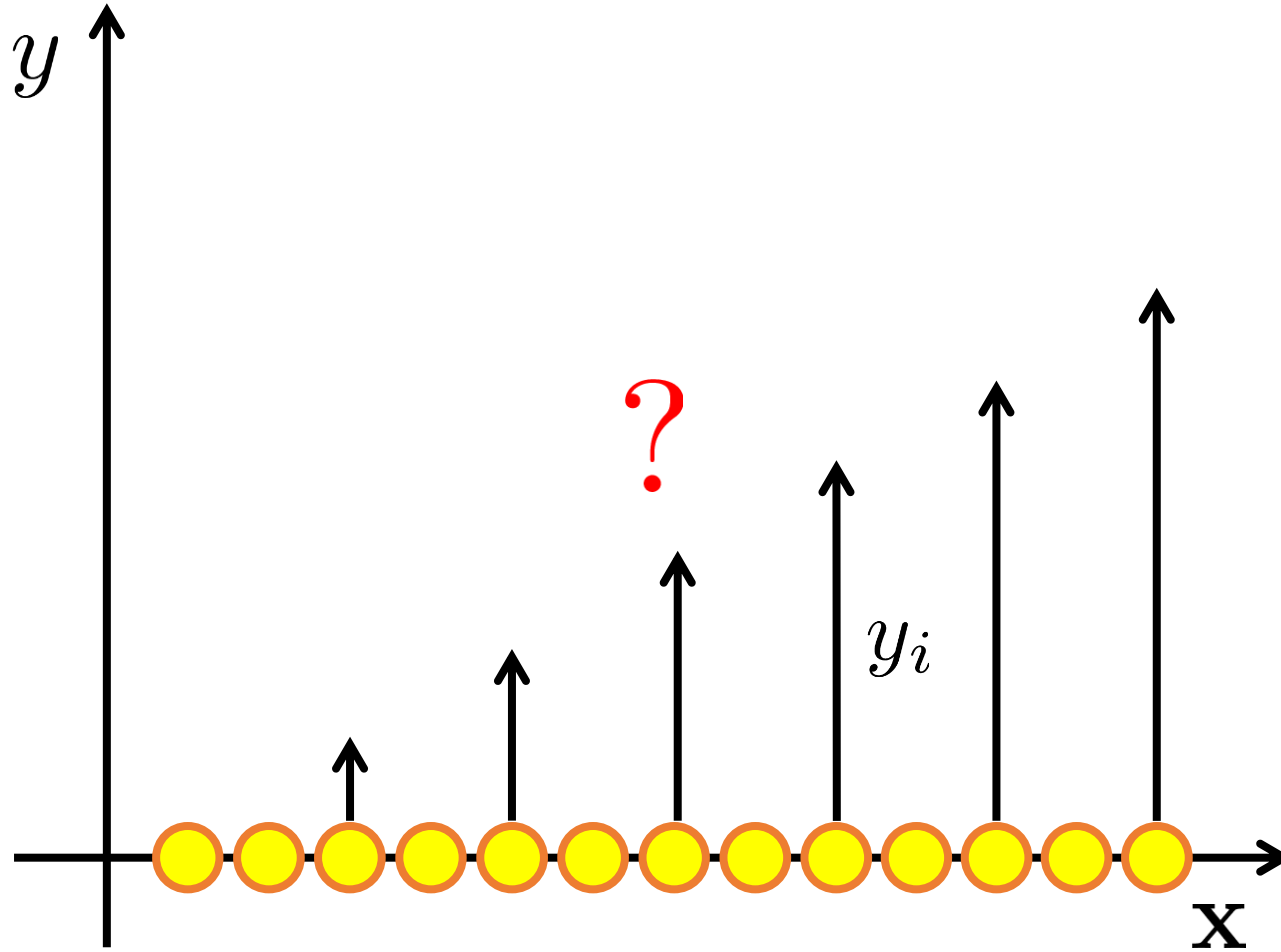


Regression using PML

August 16, 2017



Linear Regression



Data: $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n \in \mathbb{R}^d$

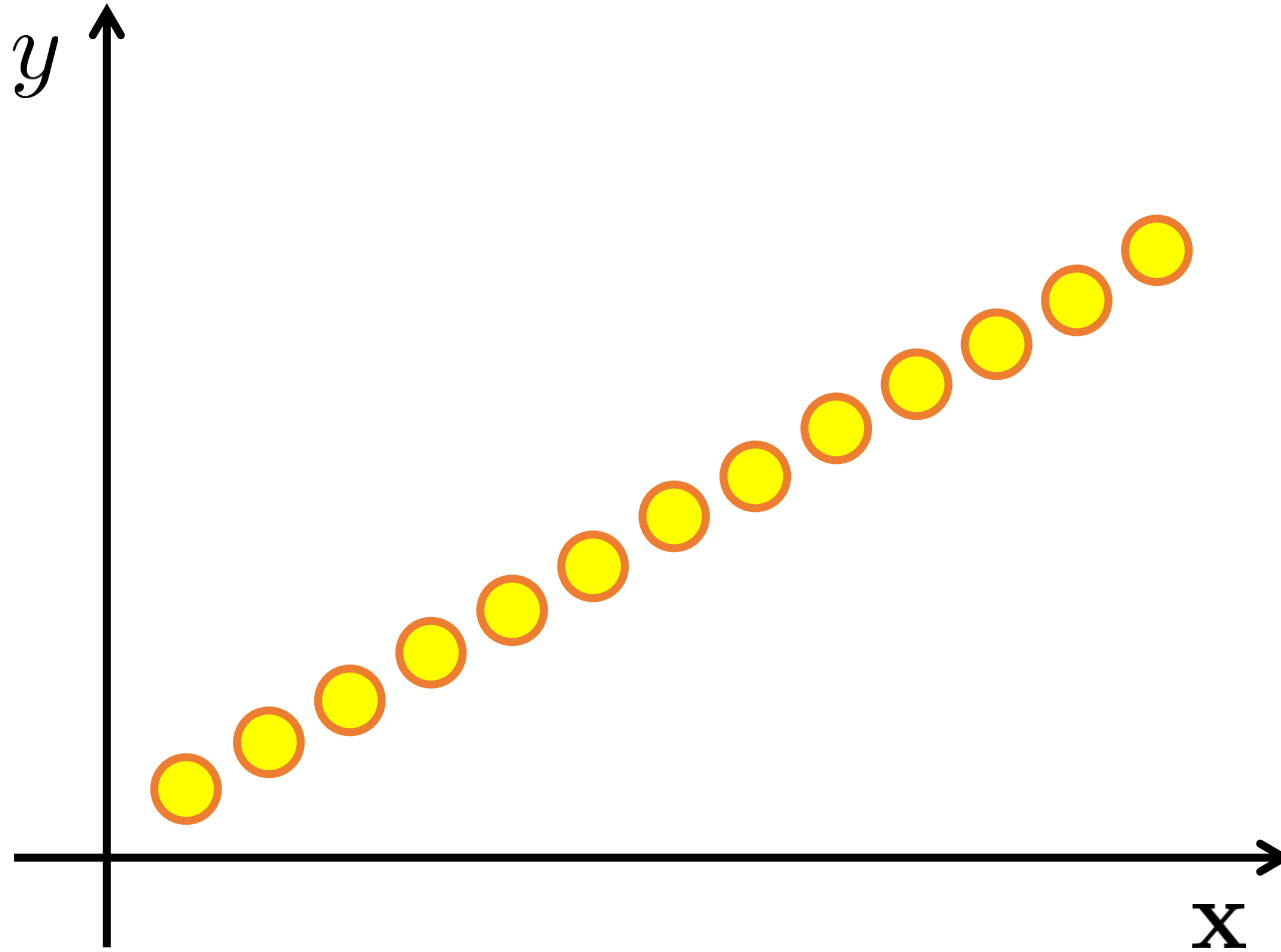
Model: \mathbf{w}^* (hidden)

$$y^i = \langle \mathbf{w}^*, \mathbf{x}^i \rangle$$

Given: $(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^n, y^n)$

Recover \mathbf{w}^* ?

Linear Regression



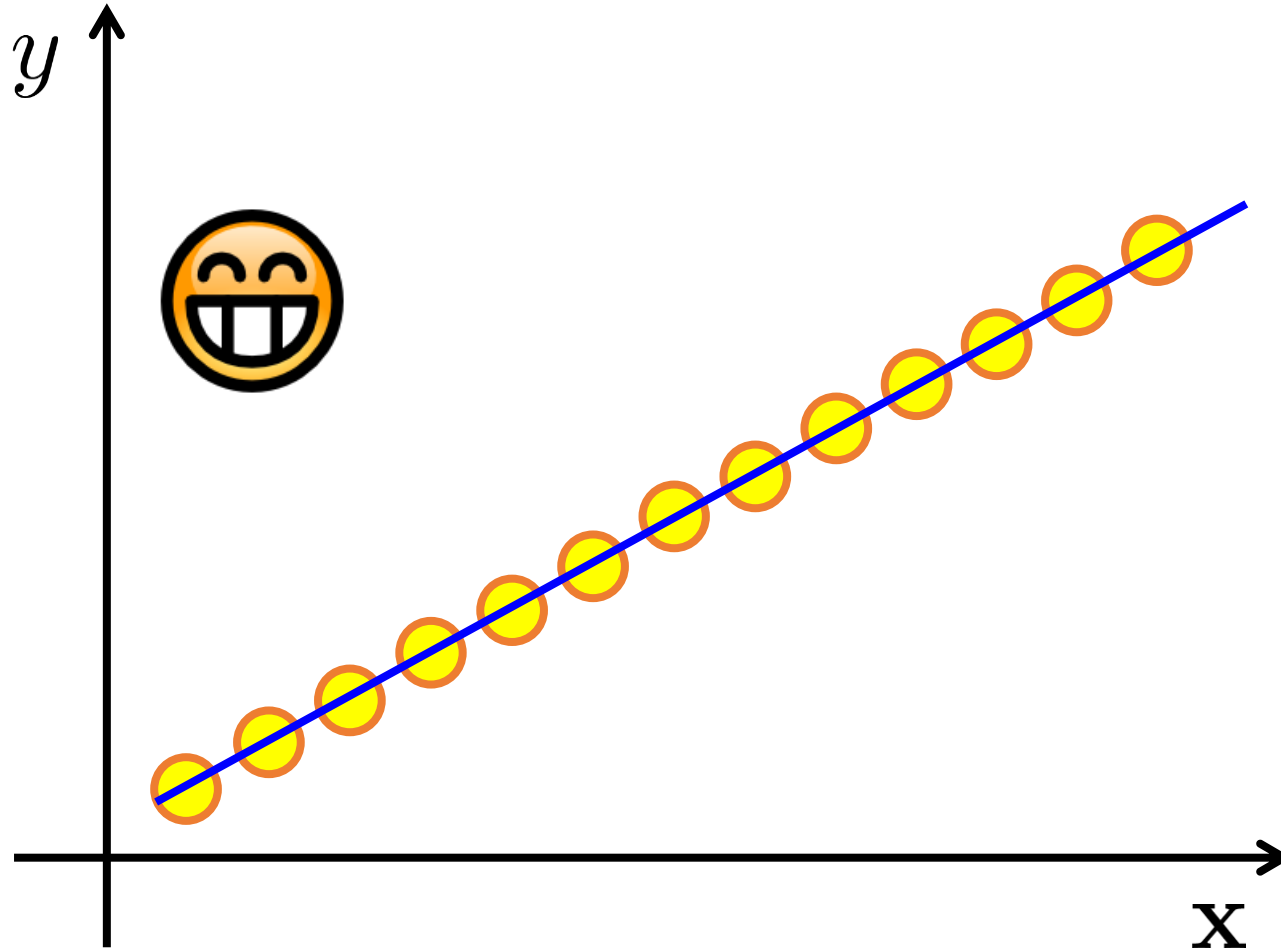
Given: $(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^n, y^n)$

$$y^i = \langle \mathbf{w}^*, \mathbf{x}^i \rangle$$

$$\begin{array}{c} \textcolor{blue}{n} \updownarrow \\ \underbrace{\begin{bmatrix} y^1 \\ y^2 \\ y^3 \\ \vdots \\ y^n \end{bmatrix}}_{\mathbf{y}} = \underbrace{\begin{bmatrix} \text{---} \mathbf{x}^1 \text{---} \\ \text{---} \mathbf{x}^2 \text{---} \\ \text{---} \mathbf{x}^3 \text{---} \\ \vdots \\ \text{---} \mathbf{x}^n \text{---} \end{bmatrix}}_{\mathbf{X}^\top} \begin{bmatrix} | \\ \mathbf{w} \\ | \end{bmatrix} \textcolor{red}{d} \updownarrow \end{array}$$

$\textcolor{red}{d}$

Linear Regression



Given: $(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^n, y^n)$

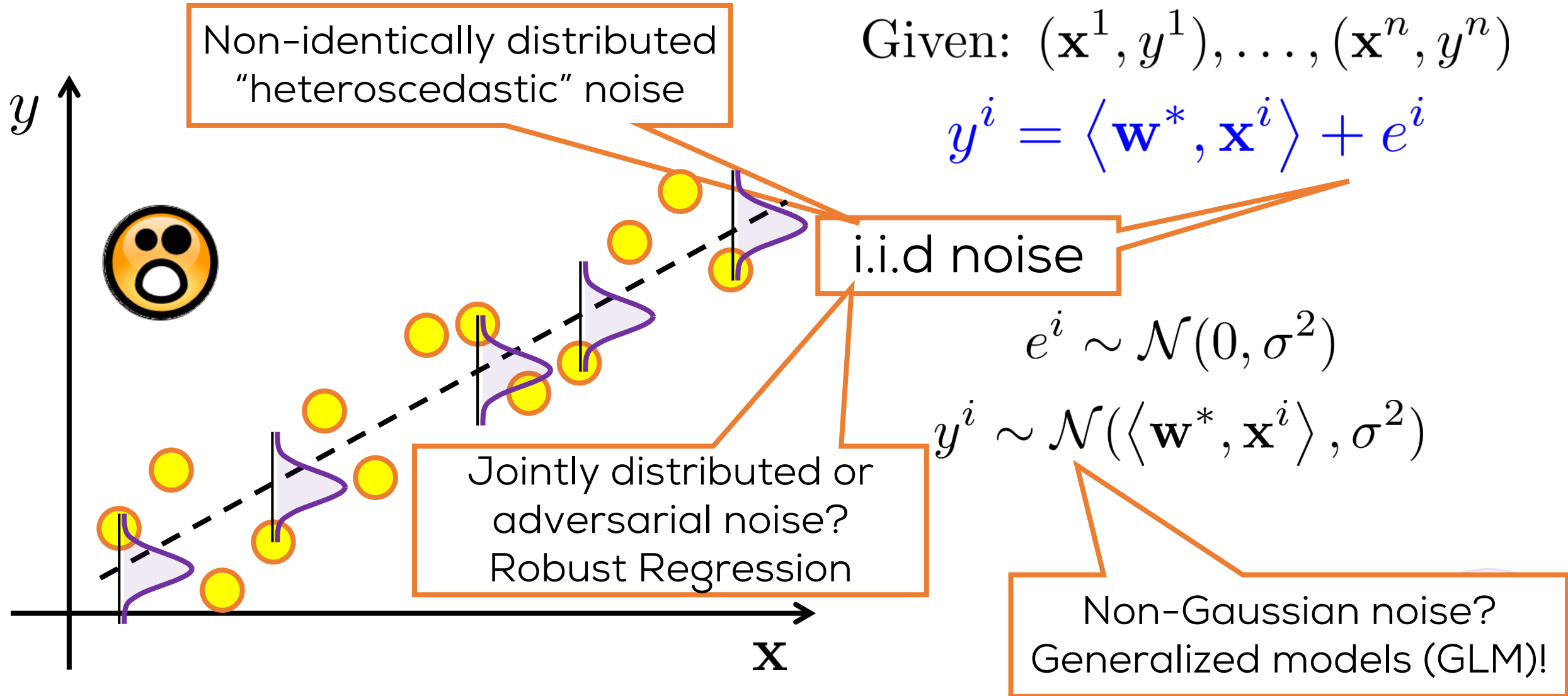
$$y^i = \langle \mathbf{w}^*, \mathbf{x}^i \rangle$$

$$\mathbf{y} = \mathbf{X}^\top \mathbf{w}$$

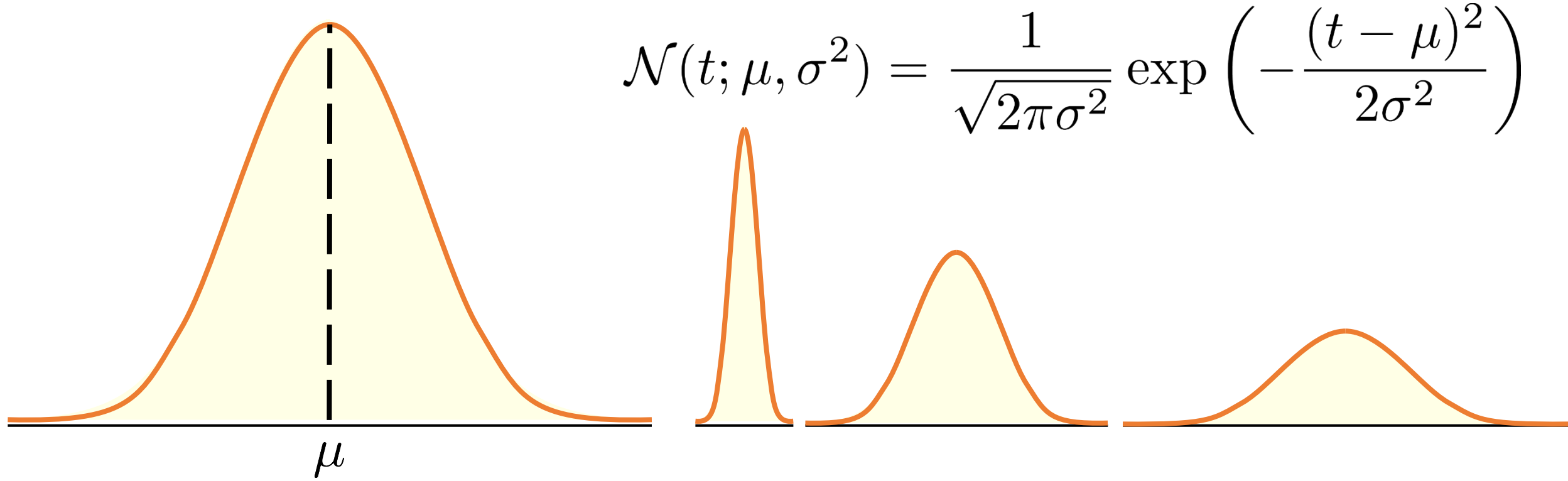
Linear system!!

\mathbf{w}^* Recovered!!

Linear Regression with Noise



The Gaussian Distribution



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

**Multivariate
Gaussian**

$$\mathcal{N}(\mathbf{v}; \boldsymbol{\mu}, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{v} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{v} - \boldsymbol{\mu})\right)$$

Linear Regression using MLE

$$\mathbb{P}[y | \mathbf{x}^i, \mathbf{w}] = \mathcal{N}(\langle \mathbf{w}, \mathbf{x}^i \rangle, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2}{2\sigma^2}\right)$$

Linear
function

Likelihood

Why isn't \mathbf{x} getting
modelled too?

Generative
model

Log-likelihood

$$\log \mathbb{P}[\mathbf{y} | \mathbf{X}, \mathbf{w}] = C - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2$$

Least Squares!

$$\hat{\mathbf{w}}_{\text{MLE}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 = (\mathbf{X}\mathbf{X}^\top)^\dagger \mathbf{X}\mathbf{y}$$

$\log \mathbb{P}[\mathbf{y}, \mathbf{X} | \boldsymbol{\theta}]$
Generative
MLE??

Exercise

Linear Regression using MAP

$$\mathbb{P}[\mathbf{w}] = \mathcal{N}(\mathbf{0}, \rho^2 \cdot I_d) = \frac{1}{\sqrt{(2\pi)^d \rho^2}} \exp\left(-\frac{\|\mathbf{w}\|_2^2}{2\rho^2}\right)$$

Sparsity inducing priors

$$\log \mathbb{P}[\mathbf{w} | \mathbf{X}, \mathbf{y}] = C - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 - \frac{1}{2\rho^2} \|\mathbf{w}\|_2^2$$

$$\hat{\mathbf{w}}_{\text{MAP}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \frac{\sigma^2}{\rho^2} \|\mathbf{w}\|_2^2 = (\mathbf{X}\mathbf{X}^\top + \lambda I)^{-1} \mathbf{X}\mathbf{y}$$

Ridge Regression

Exercise

$\log \mathbb{P}[\boldsymbol{\theta} | \mathbf{y}, \mathbf{X}]$
Generative
MAP??

Bayesian Linear Regression?

$$\mathbb{P}[\mathbf{w} \mid \mathbf{X}, \mathbf{y}] = \mathcal{N}(\boldsymbol{\nu}, \Lambda)$$

Cool! Posterior is Gaussian like the prior!

$$\boldsymbol{\nu} = \left(\mathbf{X}\mathbf{X}^\top + \frac{\sigma^2}{\rho^2} \cdot I \right)^{-1} \mathbf{X}\mathbf{y}$$
$$\Lambda = \left(\frac{1}{\sigma^2} \mathbf{X}\mathbf{X}^\top + \frac{1}{\rho^2} \cdot I \right)^{-1}$$

Wait! MAP?

By definition, MAP is the mode of the posterior

$$\mathbb{P}[y \mid \mathbf{x}, \mathbf{X}, \mathbf{y}] = \int_{\mathbf{w}} \mathbb{P}[y \mid \mathbf{x}, \mathbf{w}] \mathbb{P}[\mathbf{w} \mid \mathbf{X}, \mathbf{y}] d\mathbf{w}$$
$$= \mathcal{N}(\langle \boldsymbol{\nu}, \mathbf{x} \rangle, \sigma^2 + \mathbf{x}^\top \Lambda \mathbf{x})$$

Predictive Posterior

Extra Information

A few Thoughts

- Do I have to use these very forms for likelihood and prior?

$$\mathbb{P} [y \mid \mathbf{x}^i, \boldsymbol{\theta}] = \mathcal{D}(\langle \mathbf{w}, \mathbf{x}^i \rangle)$$

- No, however some choices make sense and make estimation easier
- Conjugacy makes life very simple if possible to achieve
- Look at your application and take these decisions
 - If you know your model is sparse, use a Laplacian prior
 - If you know your noise is not Gaussian, use a GLM
- Can I build probabilistic models for other ML tasks as well?
 - Yes, of course. We will look at classification problems next.
 - We can do probabilistic clustering, dim. redn., ranking
 - There are entire courses devoted to PML techniques: CS772, CS698S

A few Thoughts

- What about non-linear regression?

$$\mathbb{P} [y \mid \mathbf{x}^i, \boldsymbol{\theta}] = \mathcal{N}(f(\mathbf{x}^i, \boldsymbol{\theta}), \sigma^2)$$

- MLE, MAP estimation more challenging - kernel methods, deep learning
- Gaussian processes more suited for non-linear PML – beyond the scope!
- Are Bayesian models better than non-Bayesian ones?
 - Bayesian models more informative $\mathcal{N}(\langle \boldsymbol{\nu}, \mathbf{x} \rangle, \sigma^2 + \mathbf{x}^\top \Lambda \mathbf{x})$
 - Useful in settings like active learning, anomaly detection
 - Can be more expensive at training and prediction time
 - Ask your doctor if your application needs Bayesian reasoning or not!
- Bayesian \neq Generative

A few Thoughts

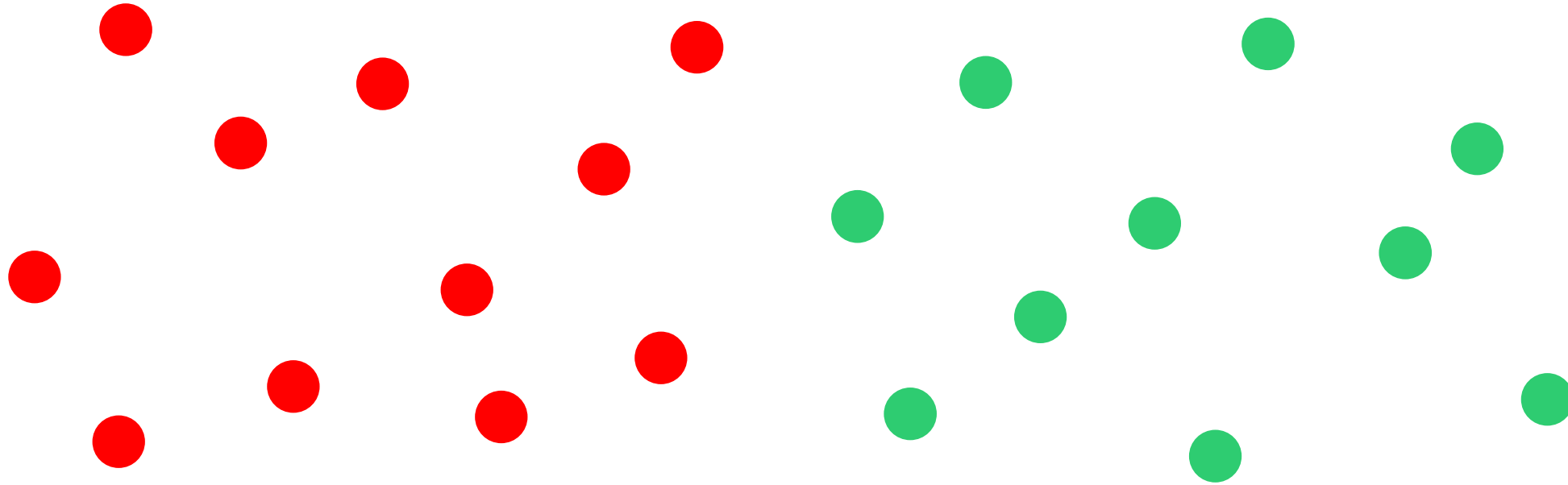
- How do I solve all these optimization problems that MLE, MAP procedures keep throwing up?
 - We are in safe hands – optimization theory is very evolved
 - Wait for the next set of lectures on function approximation methods
- Are generatively trained models better or discriminative ones?
 - Discriminative models reason about $\mathbb{P}[y|\mathbf{x}]$
 - Generative models reason about $\mathbb{P}[y, \mathbf{x}]$
 - For prediction, all we need is $\mathbb{P}[y|\mathbf{x}]$
 - Generative models handle missing, erroneous data easily but bulky
 - Discriminative models are much lighter, easier to train
 - Models such GMMs, deep nets, can be trained in gen. and disc. manner

Classification using PML

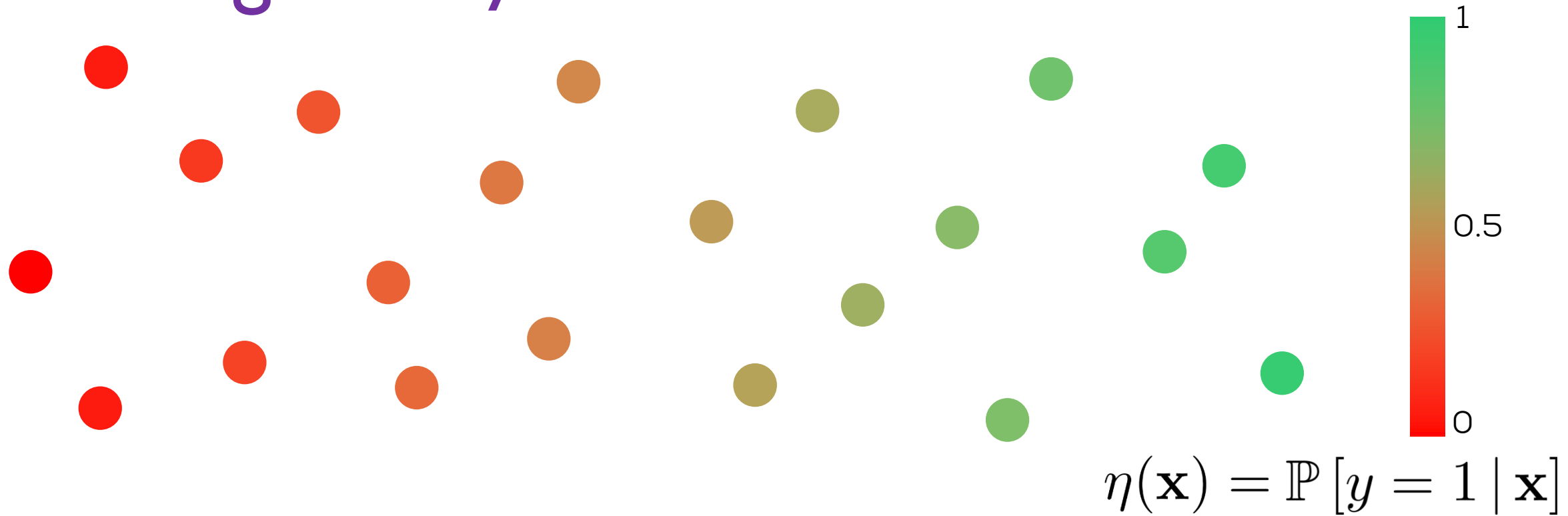
August 16, 2017



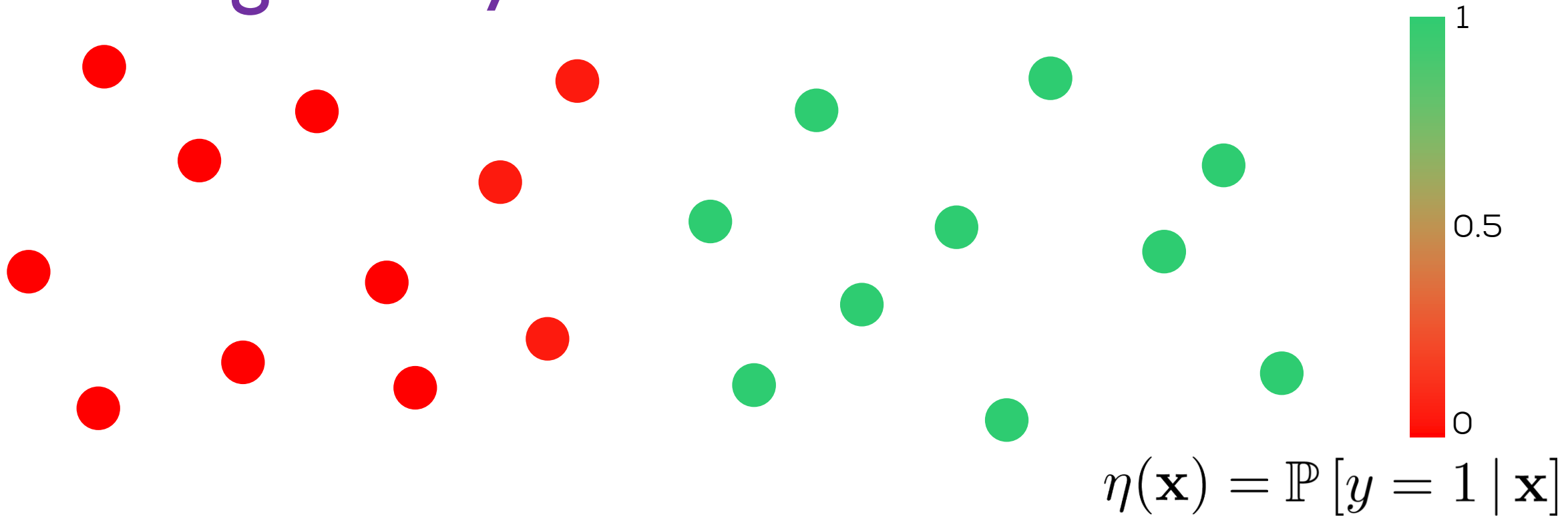
Revisiting Binary Classification



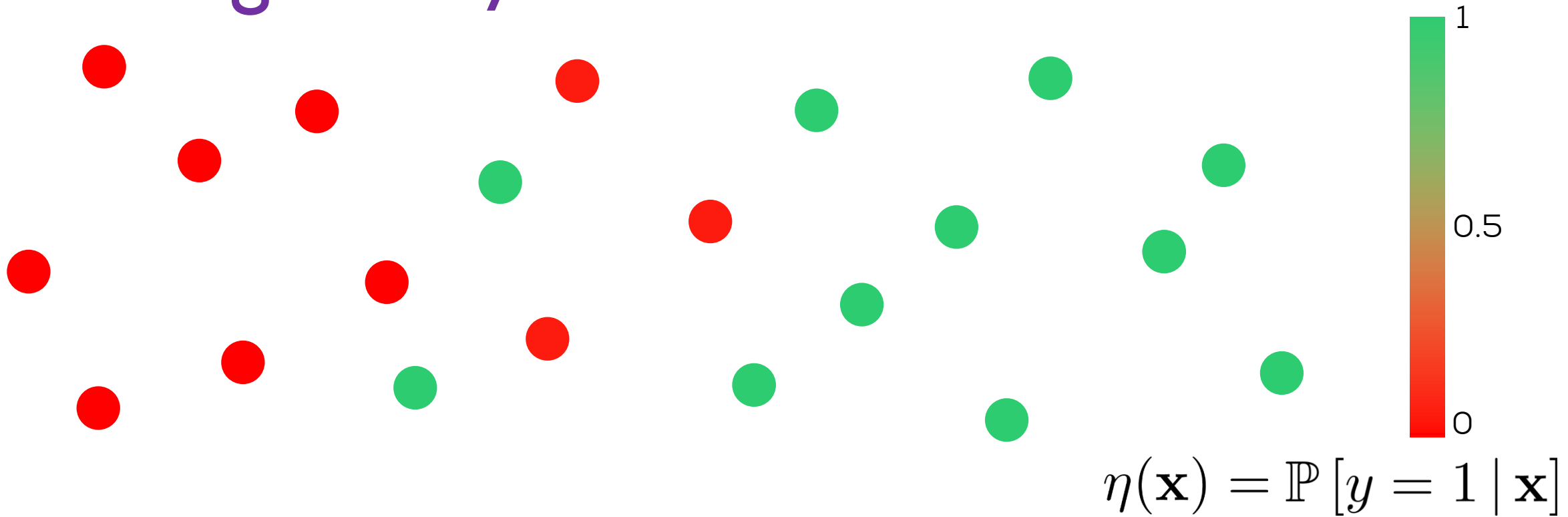
Revisiting Binary Classification



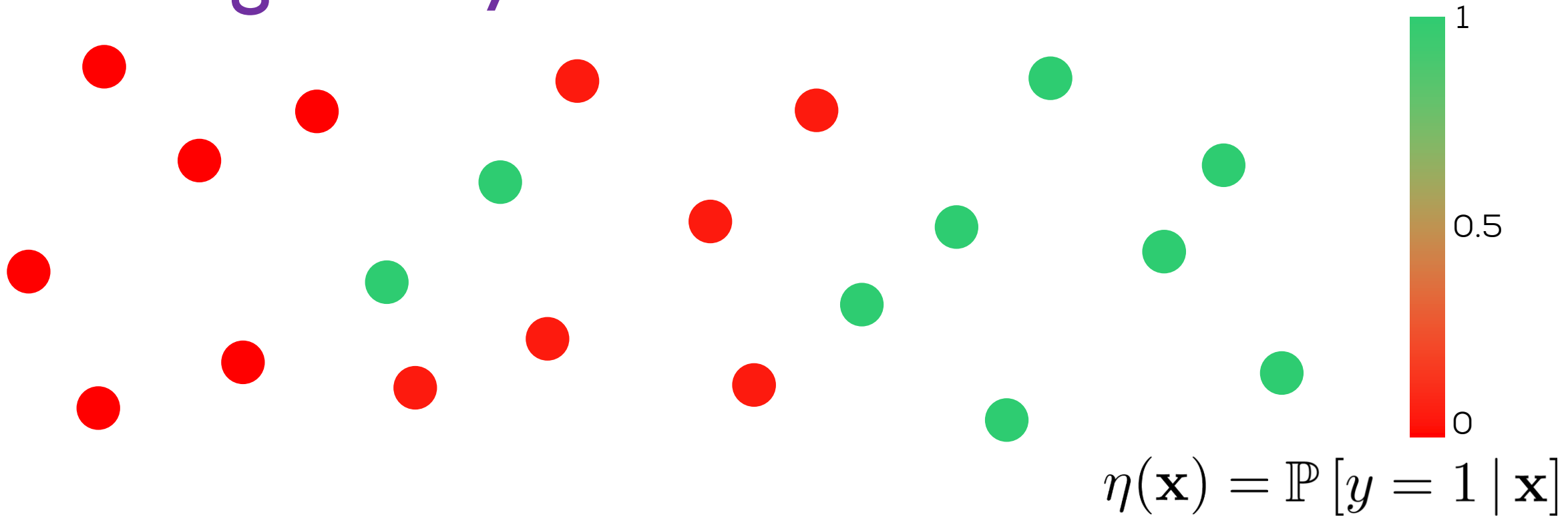
Revisiting Binary Classification



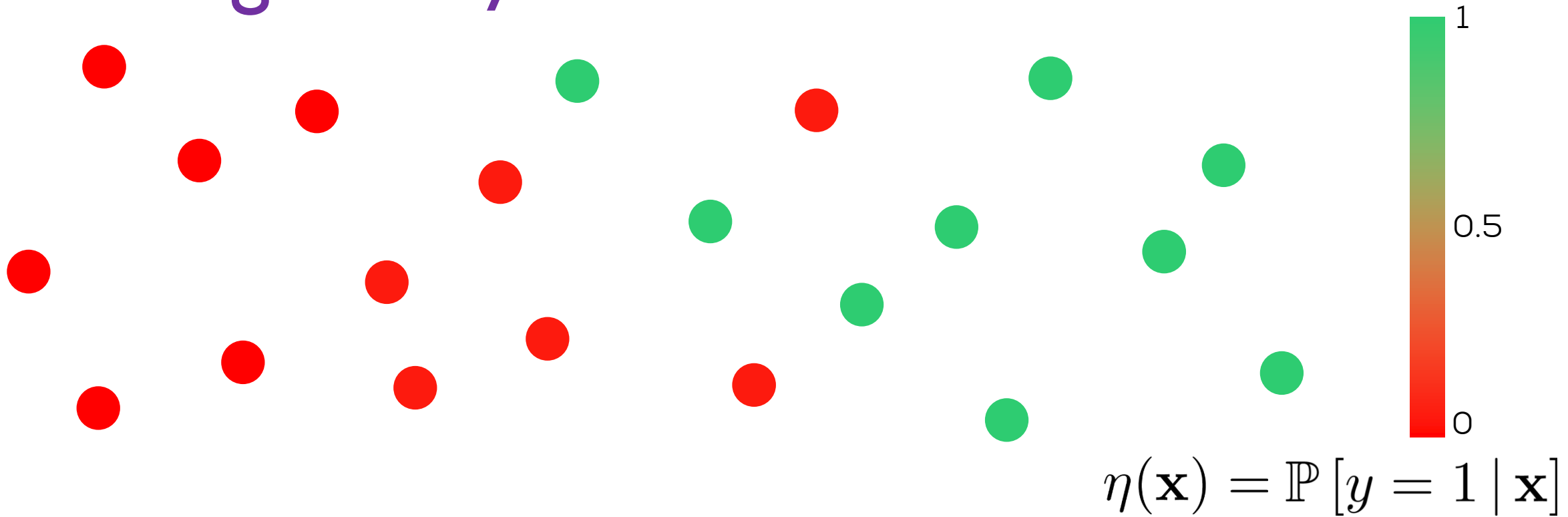
Revisiting Binary Classification



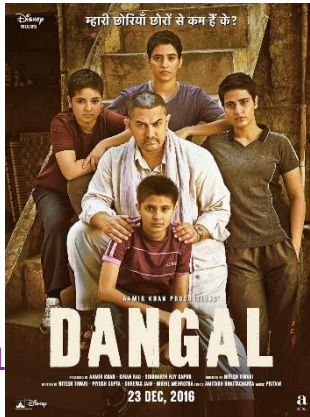
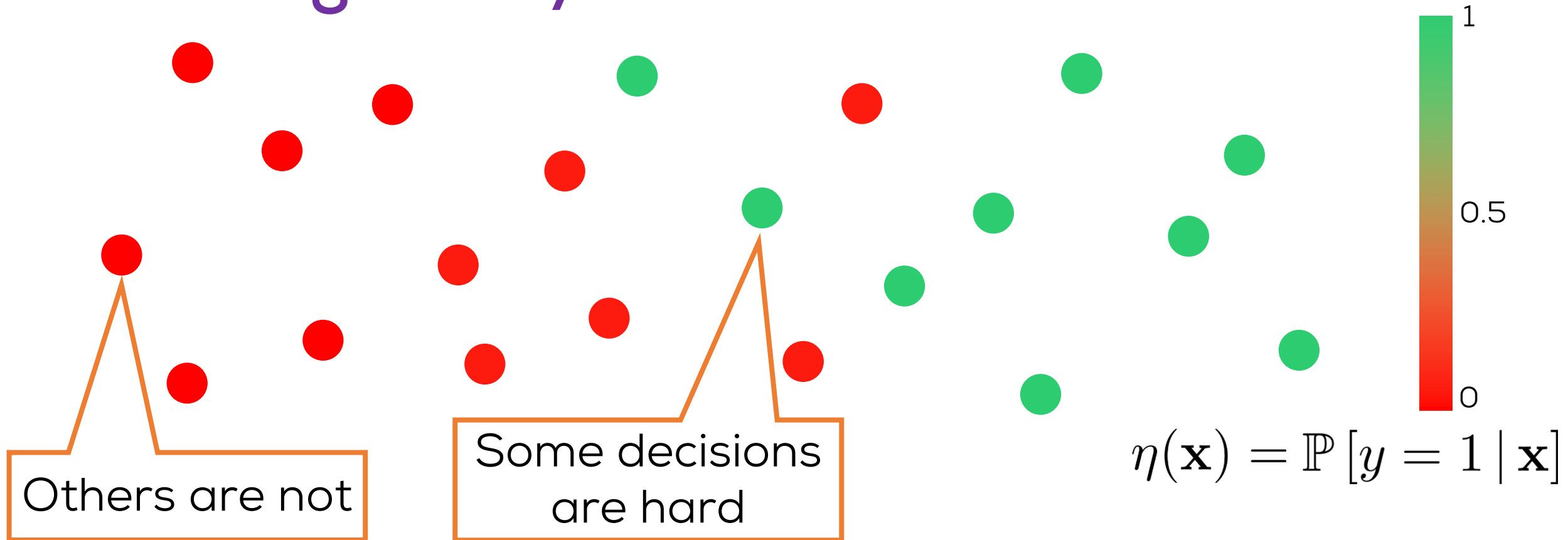
Revisiting Binary Classification



Revisiting Binary Classification

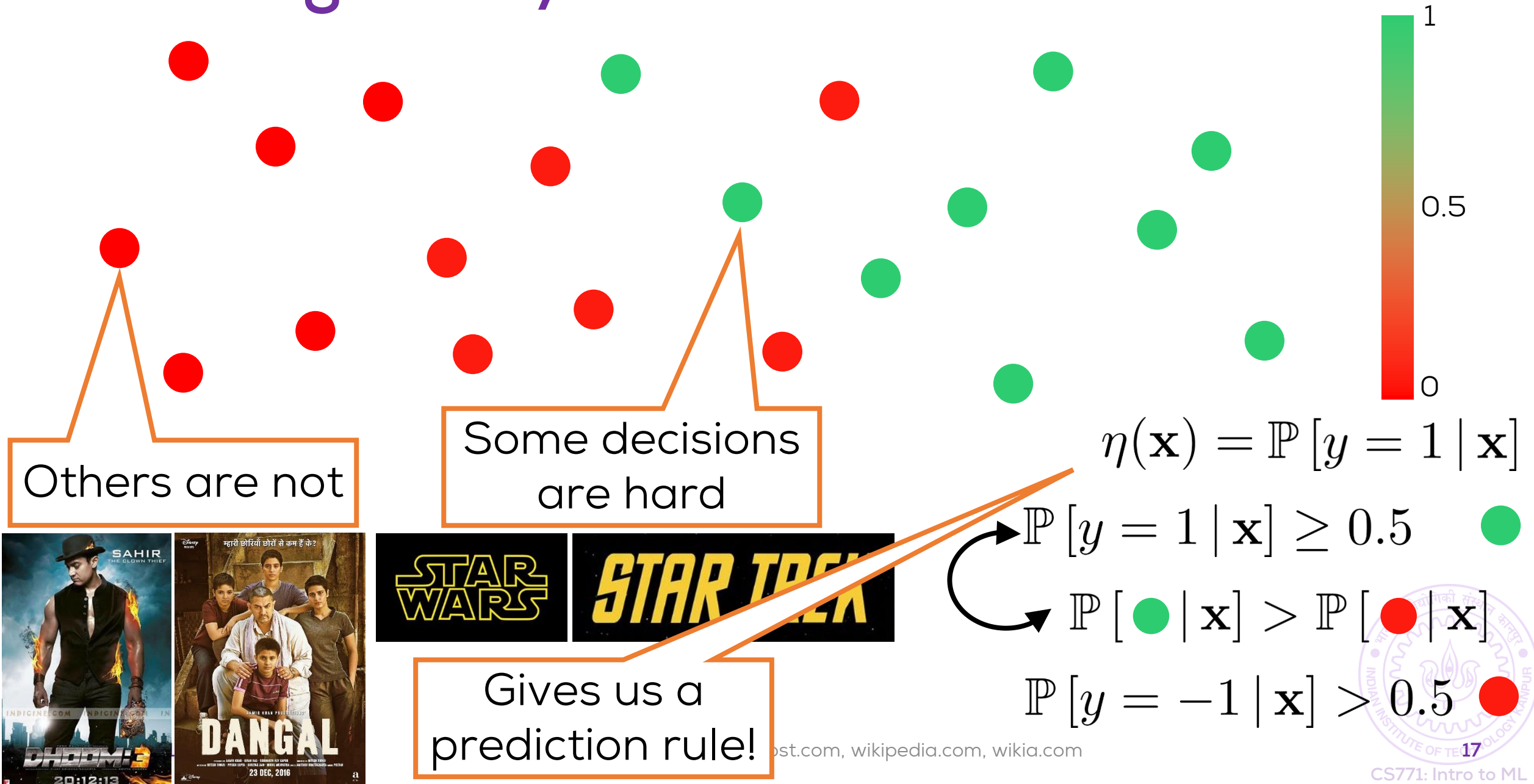


Revisiting Binary Classification

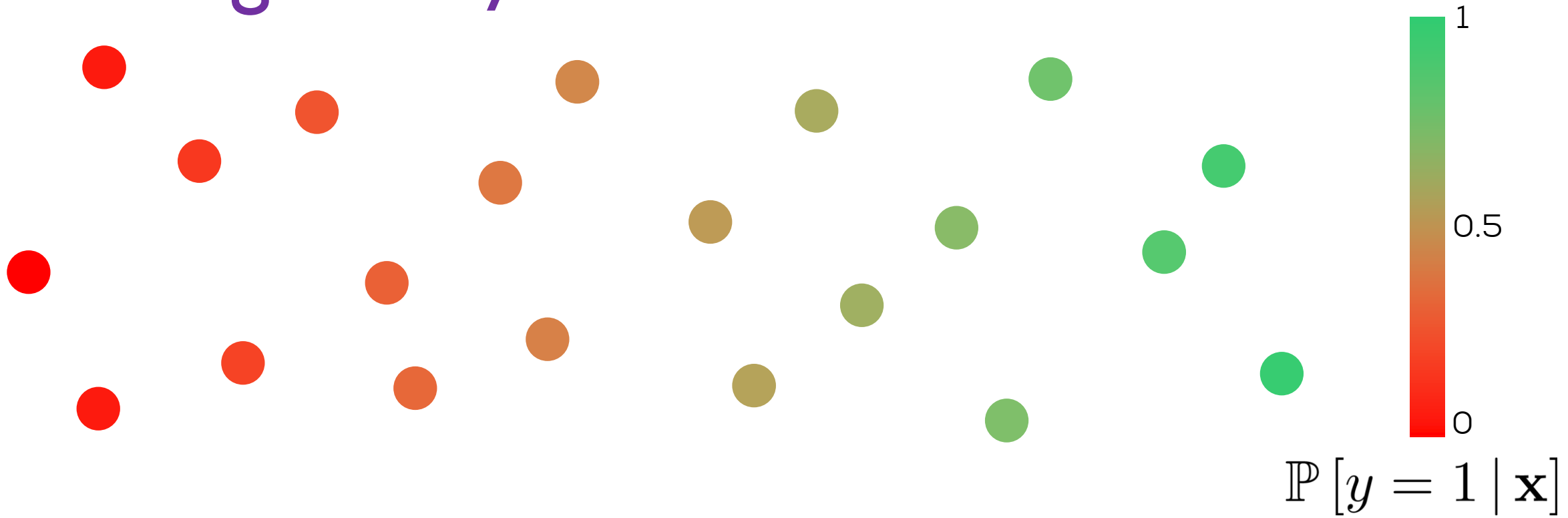


subscene.com, firstpost.com, wikipedia.com, wikia.com

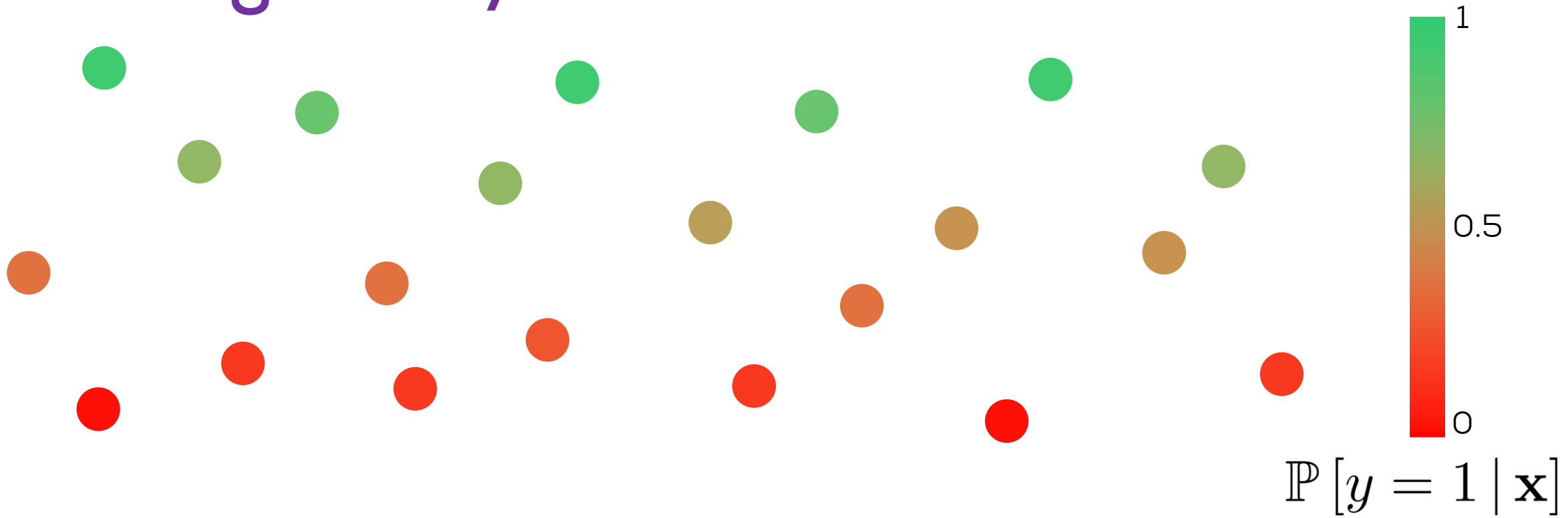
Revisiting Binary Classification



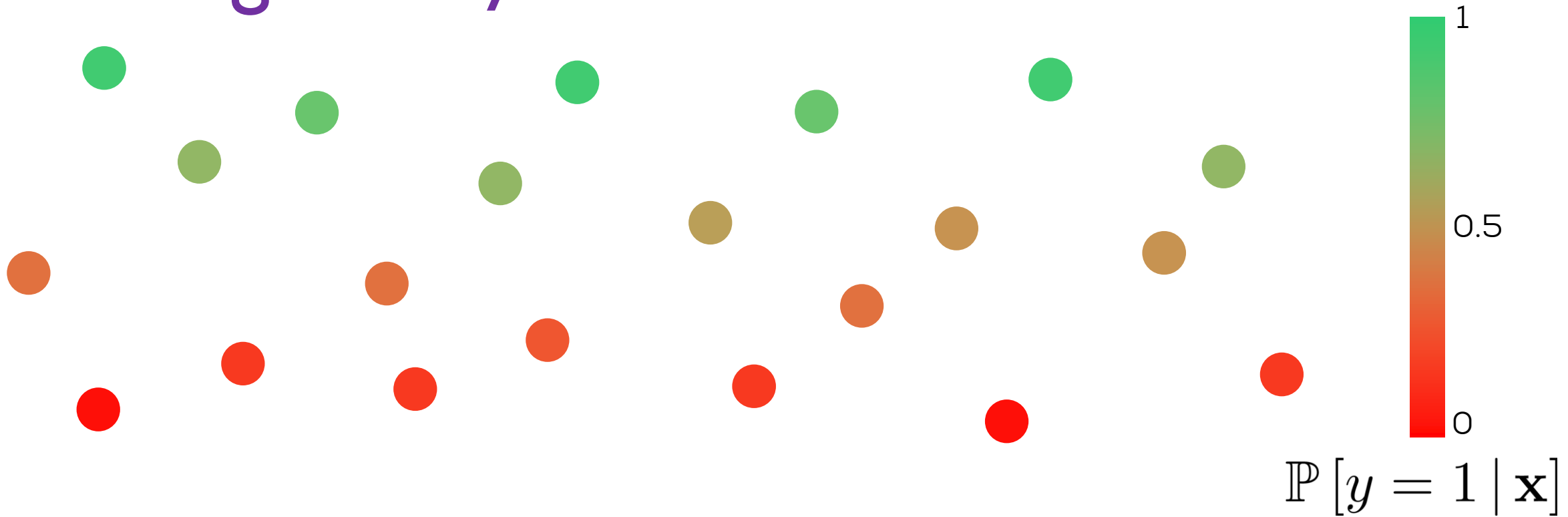
Revisiting Binary Classification



Revisiting Binary Classification

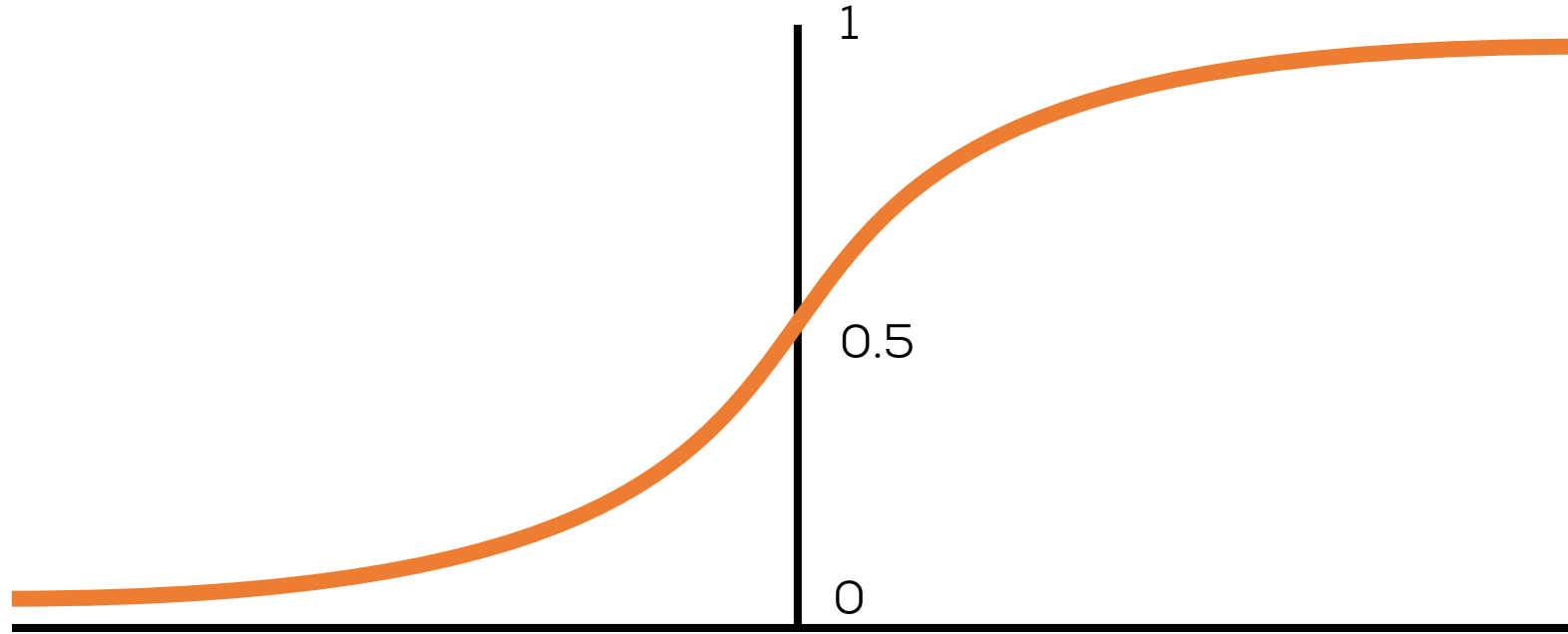


Revisiting Binary Classification



Features Matter!

The Sigmoid Function



$$\sigma(t) = \frac{1}{1 + \exp(-t)} = \frac{\exp(t)}{\exp(t) + 1}$$

Classification using MLE

$$\mathbb{P}[y | \mathbf{x}^i, \mathbf{w}] = \sigma(y \langle \mathbf{w}, \mathbf{x}^i \rangle) = \frac{1}{1 + \exp(-y \langle \mathbf{w}, \mathbf{x}^i \rangle)}$$

Shorthand
 $\mathbb{P}[y^i = y] =: \mathbb{P}[y]$

$$\mathbb{P}[y^i = 1 | \mathbf{x}^i, \mathbf{w}] + \mathbb{P}[y^i = -1 | \mathbf{x}^i, \mathbf{w}] = 1$$

Likelihood

$$\log \frac{\mathbb{P}[y^i = 1 | \mathbf{x}^i, \mathbf{w}]}{\mathbb{P}[y^i = -1 | \mathbf{x}^i, \mathbf{w}]} = \langle \mathbf{w}, \mathbf{x}^i \rangle$$

Log-odds

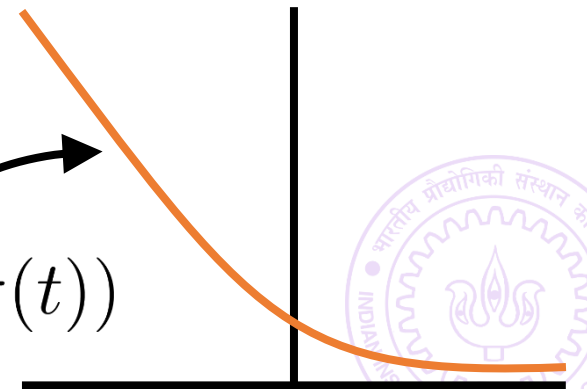
Log-likelihood

$$\log \mathbb{P}[\mathbf{y} | \mathbf{X}, \mathbf{w}] = \sum_{i=1}^n \log \left(\frac{1}{1 + \exp(-y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)} \right)$$

$\hat{\mathbf{w}}_{\text{MLE}} = ?$

Wait for learning
with FA

$$f(t) = -\log(\sigma(t))$$



Another possibility

- What happens if $y_i \in \{0,1\}$ instead of $y_i \in \{-1,1\}$?
- Binomial instead of Rademacher
- How to redefine likelihood?
- Need to ensure $\mathbb{P}[y^i = 1 \mid \mathbf{x}^i, \mathbf{w}] + \mathbb{P}[y^i = 0 \mid \mathbf{x}^i, \mathbf{w}] = 1$
- One way is the following

$$\hat{\eta}_i = \sigma(\langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\mathbb{P}[y^i \mid \mathbf{x}^i, \mathbf{w}] = (\hat{\eta}_i)^{y^i} \cdot (1 - \hat{\eta}_i)^{1-y^i}$$

- Do you get the same MLE problem?
- What about MAP estimates? $\mathbb{P}[\mathbf{w}] = \mathcal{N}(\mathbf{0}, \rho^2 \cdot I_d)$

Bayesian Logistic Regression?



$$\mathbb{P}[y \mid \mathbf{x}, \mathbf{w}] = \sigma(y \langle \mathbf{w}, \mathbf{x} \rangle)$$

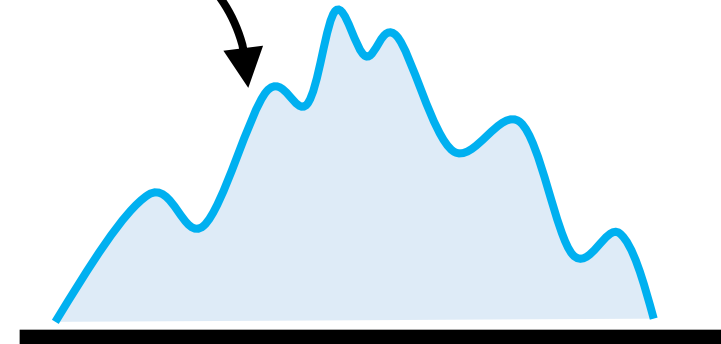
$$\mathbb{P}[\mathbf{w}] = \mathcal{N}(\mathbf{0}, \rho^2 \cdot I_d)$$

$$\mathbb{P}[\mathbf{w} \mid \mathbf{X}, \mathbf{y}] = \frac{\mathbb{P}[\mathbf{y} \mid \mathbf{X}, \mathbf{w}] \mathbb{P}[\mathbf{w}]}{\int_{\mathbf{w}'} \mathbb{P}[\mathbf{y} \mid \mathbf{X}, \mathbf{w}'] \mathbb{P}[\mathbf{w}'] d\mathbf{w}'} = \frac{\prod_{i=1}^n \sigma(y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)}{\int_{\mathbf{w}'} \prod_{i=1}^n \sigma(y^i \langle \mathbf{w}', \mathbf{x}^i \rangle) \exp\left(-\frac{\|\mathbf{w}'\|_2^2}{2\sigma^2}\right) d\mathbf{w}'}$$

Posterior Approximation

MCMC, Variational Inference,
Laplace approx

$\mathbb{P}[\mathbf{w} \mid \mathbf{X}, \mathbf{y}]$



Bayesian Logistic Regression?



$$\mathbb{P}[y | \mathbf{x}, \mathbf{w}] = \sigma(y \langle \mathbf{w}, \mathbf{x} \rangle)$$

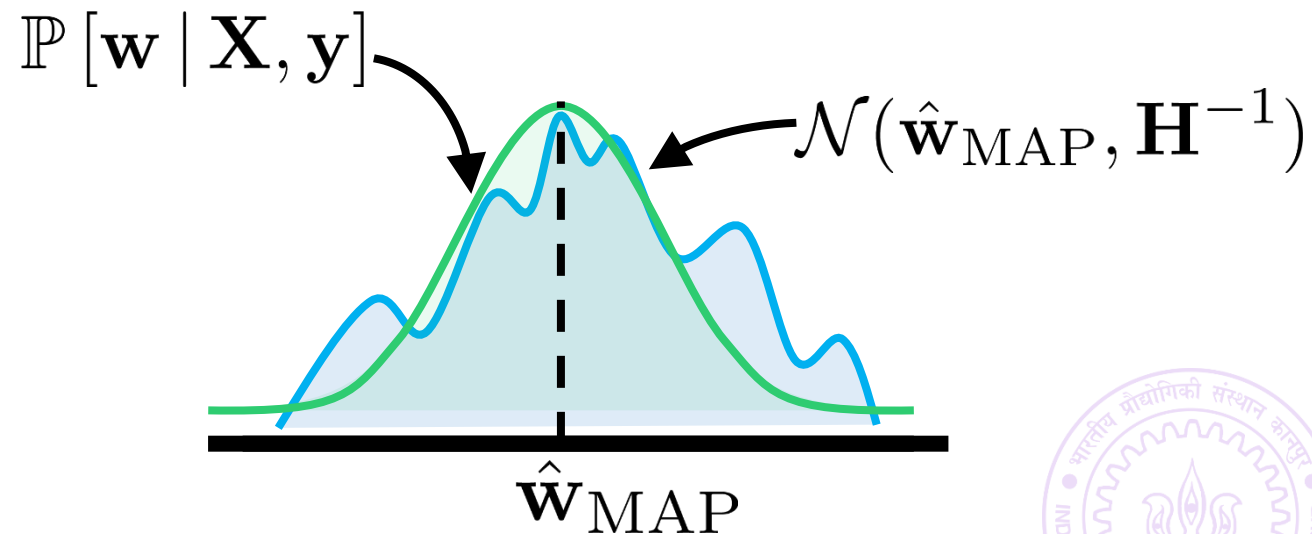
$$\mathbb{P}[\mathbf{w}] = \mathcal{N}(\mathbf{0}, \rho^2 \cdot I_d)$$

$$\mathbb{P}[\mathbf{w} | \mathbf{X}, \mathbf{y}] = \frac{\mathbb{P}[\mathbf{y} | \mathbf{X}, \mathbf{w}] \mathbb{P}[\mathbf{w}]}{\int_{\mathbf{w}'} \mathbb{P}[\mathbf{y} | \mathbf{X}, \mathbf{w}'] \mathbb{P}[\mathbf{w}'] d\mathbf{w}'} = \frac{\prod_{i=1}^n \sigma(y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)}{\int_{\mathbf{w}'} \prod_{i=1}^n \sigma(y^i \langle \mathbf{w}', \mathbf{x}^i \rangle) \exp\left(-\frac{\|\mathbf{w}'\|_2^2}{2\sigma^2}\right) d\mathbf{w}'}$$

Posterior Approximation

MCMC, Variational Inference,
Laplace approx

However, even the Laplace predictive posterior is intractable ☹



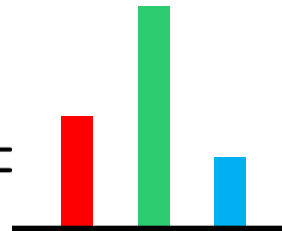
Multi-class and Multi-label Classification using PML

August 16, 2017



Multi-classification using MLE

- $K > 2$ classes – need more detailed parameters
- For each point, its label profile is a vector

$$\boldsymbol{\eta}(\mathbf{x}) =$$


$$\mathbb{P} [y^i = k \mid \mathbf{x}^i, \{\mathbf{w}^l\}_{1,\dots,K}] \propto \exp(\langle \mathbf{w}^k, \mathbf{x}^i \rangle)$$

$$\mathbb{P} [y^i = k \mid \mathbf{x}^i, \{\mathbf{w}^l\}_{1,\dots,K}] = \frac{\exp(\langle \mathbf{w}^k, \mathbf{x}^i \rangle)}{\sum_{l=1}^K \exp(\langle \mathbf{w}^l, \mathbf{x}^i \rangle)}$$

- Likelihood function is multinomial instead of binomial

$$\mathbb{P} [\mathbf{y} \mid \mathbf{X}, \mathbf{w}] = \prod_{i=1}^n \hat{\eta}_{y^i}^i(\mathbf{x}) \quad \hat{\eta}_k^i(\mathbf{x}) = \frac{\exp(\langle \mathbf{w}^k, \mathbf{x}^i \rangle)}{\sum_{l=1}^K \exp(\langle \mathbf{w}^l, \mathbf{x}^i \rangle)}$$

Softmax Regression

Binary vs Multi-Classification using Linear LR

- Binary

$$\arg \max_{y \in \{-1, 1\}} \mathbb{P}[y \mid \mathbf{x}]$$

- Multiclass

Binary vs Multi-Classification using Linear LR

- Binary

$$\mathbb{P}[y \mid \mathbf{x}] \geq 0.5$$

- Multiclass



Binary vs Multi-Classification using Linear LR

- Binary

$$\mathbb{P}[y \mid \mathbf{x}] \geq 0.5$$

$$\mathbb{P}[y \mid \mathbf{x}, \mathbf{w}] = \sigma(y \langle \mathbf{w}, \mathbf{x} \rangle)$$

- Multiclass



Binary vs Multi-Classification using Linear LR

- Binary

$$\sigma(y \langle \mathbf{w}, \mathbf{x}^i \rangle) \geq 0.5$$

- Multiclass



Binary vs Multi-Classification using Linear LR

- Binary

$$y \langle \mathbf{w}, \mathbf{x}^i \rangle \geq 0$$

- Multiclass

Binary vs Multi-Classification using Linear LR

- Binary

$$y = \text{sign}(\langle \mathbf{w}, \mathbf{x}^i \rangle)$$

- Multiclass

Binary vs Multi-Classification using Linear LR

- Binary

$$y = \text{sign}(\langle \mathbf{w}, \mathbf{x}^i \rangle)$$

- Multiclass

Binary vs Multi-Classification using Linear LR

- Binary

$$y = \text{sign}(\langle \mathbf{w}, \mathbf{x}^i \rangle)$$

- Multiclass

$$\arg \max_{y \in [K]} \mathbb{P}[y | \mathbf{x}]$$

Binary vs Multi-Classification using Linear LR

- Binary

$$y = \text{sign}(\langle \mathbf{w}, \mathbf{x}^i \rangle)$$

- Multiclass

$$\arg \max_{y \in [K]} \mathbb{P}[y | \mathbf{x}]$$

$$\mathbb{P}[y = k | \mathbf{x}, \mathbf{w}] = \sigma(\langle \mathbf{w}^k, \mathbf{x} \rangle)$$

Binary vs Multi-Classification using Linear LR

- Binary

$$y = \text{sign}(\langle \mathbf{w}, \mathbf{x}^i \rangle)$$

- Multiclass

$$\arg \max_{y \in [K]} \sigma(\langle \mathbf{w}^k, \mathbf{x} \rangle)$$

Binary vs Multi-Classification using Linear LR

- Binary

$$y = \text{sign}(\langle \mathbf{w}, \mathbf{x}^i \rangle)$$

- Multiclass

$$\arg \max_{y \in [K]} \langle \mathbf{w}^k, \mathbf{x} \rangle$$

Binary vs Multi-Classification using Linear LR

- Binary

$$y = \text{sign}(\langle \mathbf{w}, \mathbf{x}^i \rangle)$$

- Multiclass

$$\arg \max_{y \in [K]} \langle \mathbf{w}^k, \mathbf{x} \rangle$$

Binary vs Multi-Classification using Linear LR

- Binary

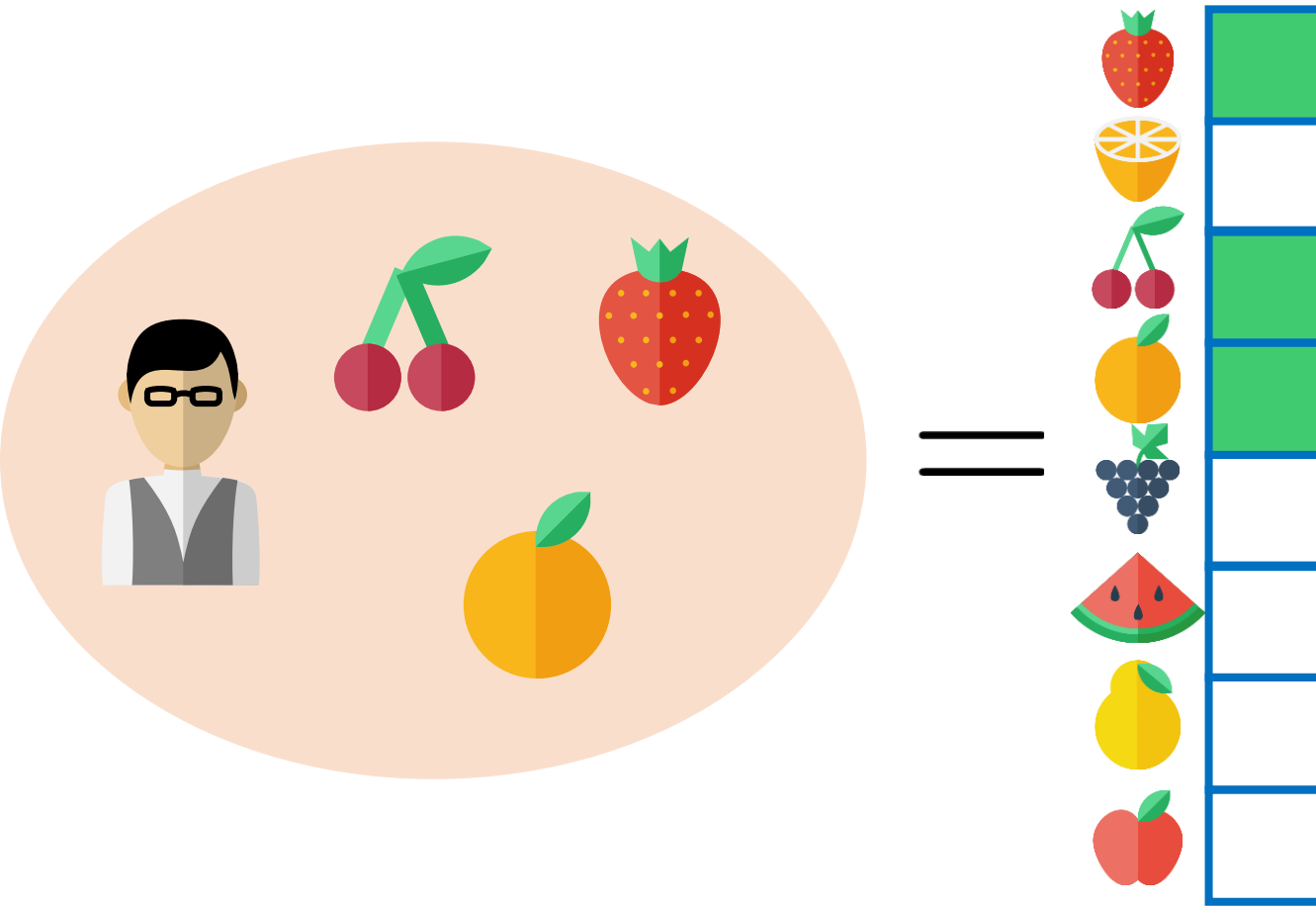
$$y = \text{sign}(\langle \mathbf{w}, \mathbf{x}^i \rangle)$$

- Multiclass

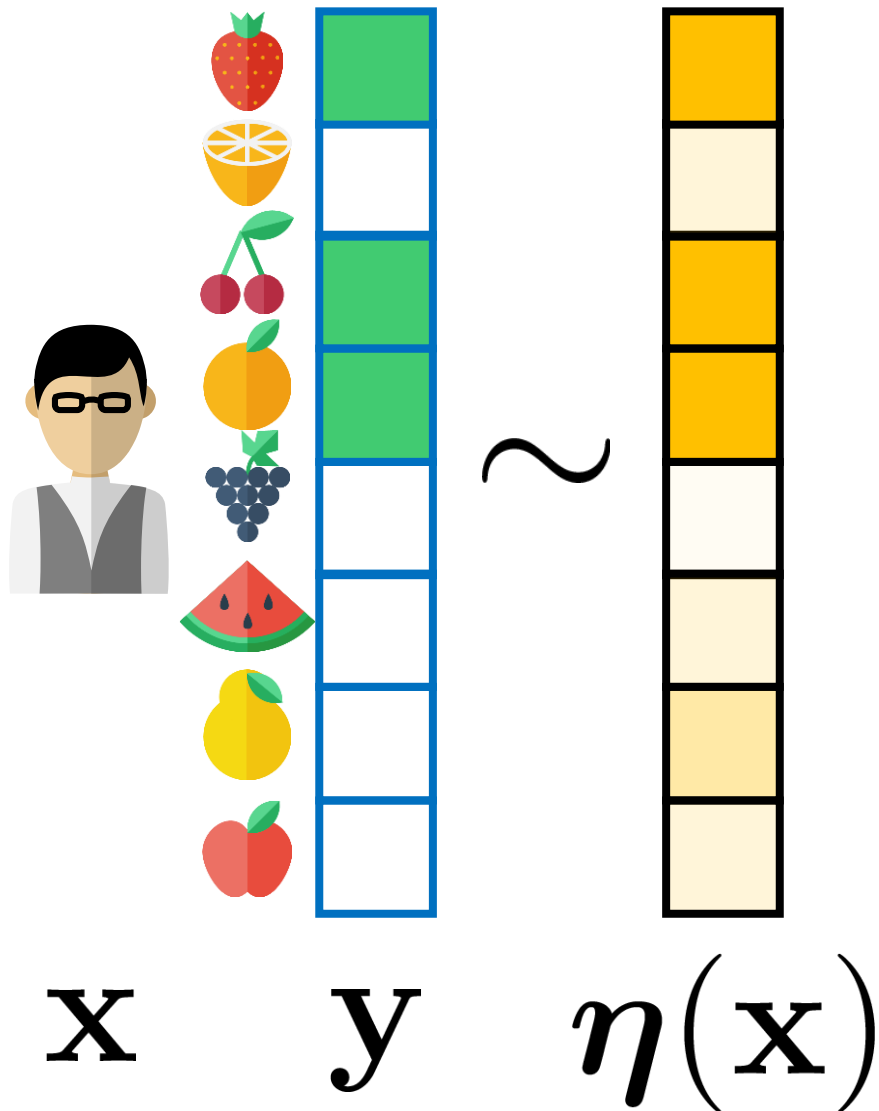
$$\arg \max_{y \in [K]} \langle \mathbf{w}^k, \mathbf{x} \rangle$$

Linear Classifier

Multi-label Learning using PML



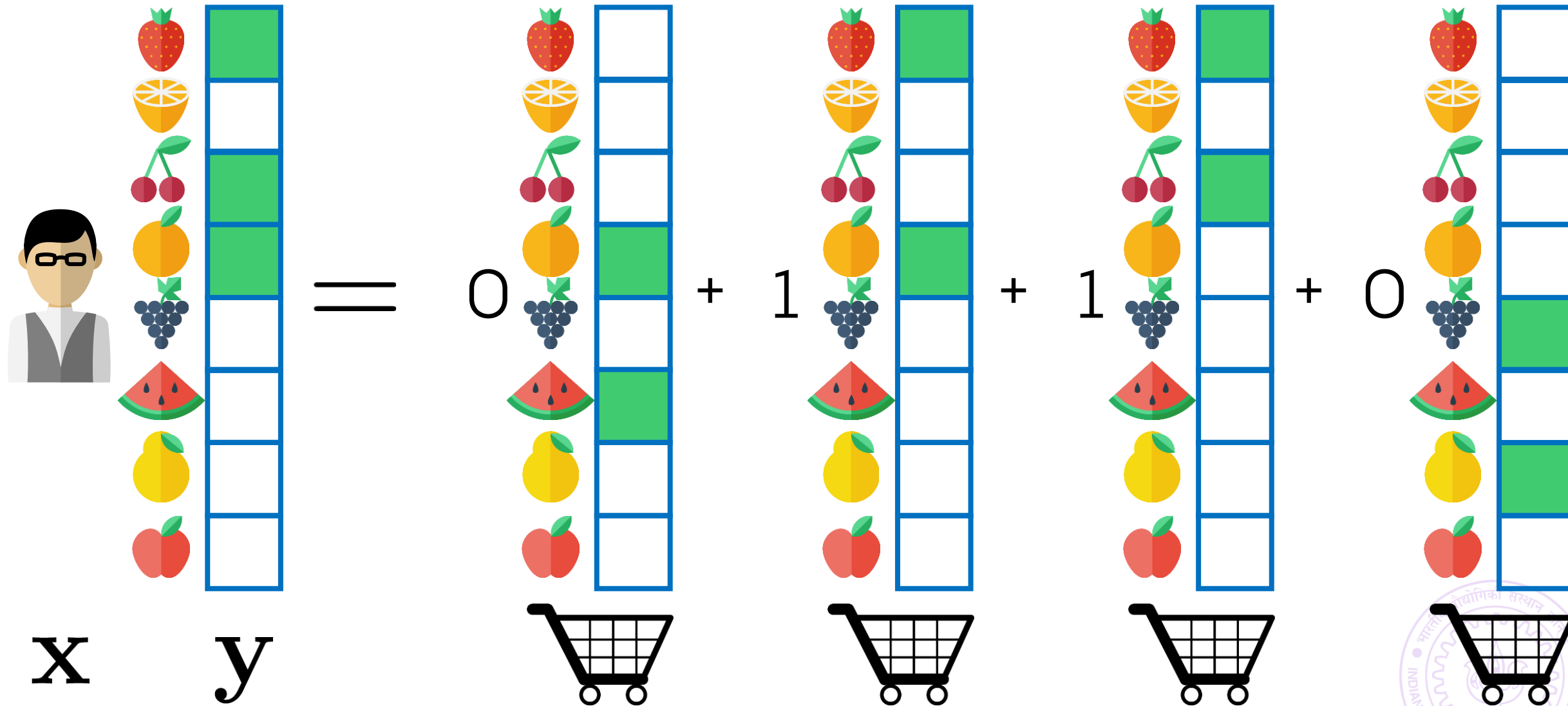
Multi-label Learning using PML



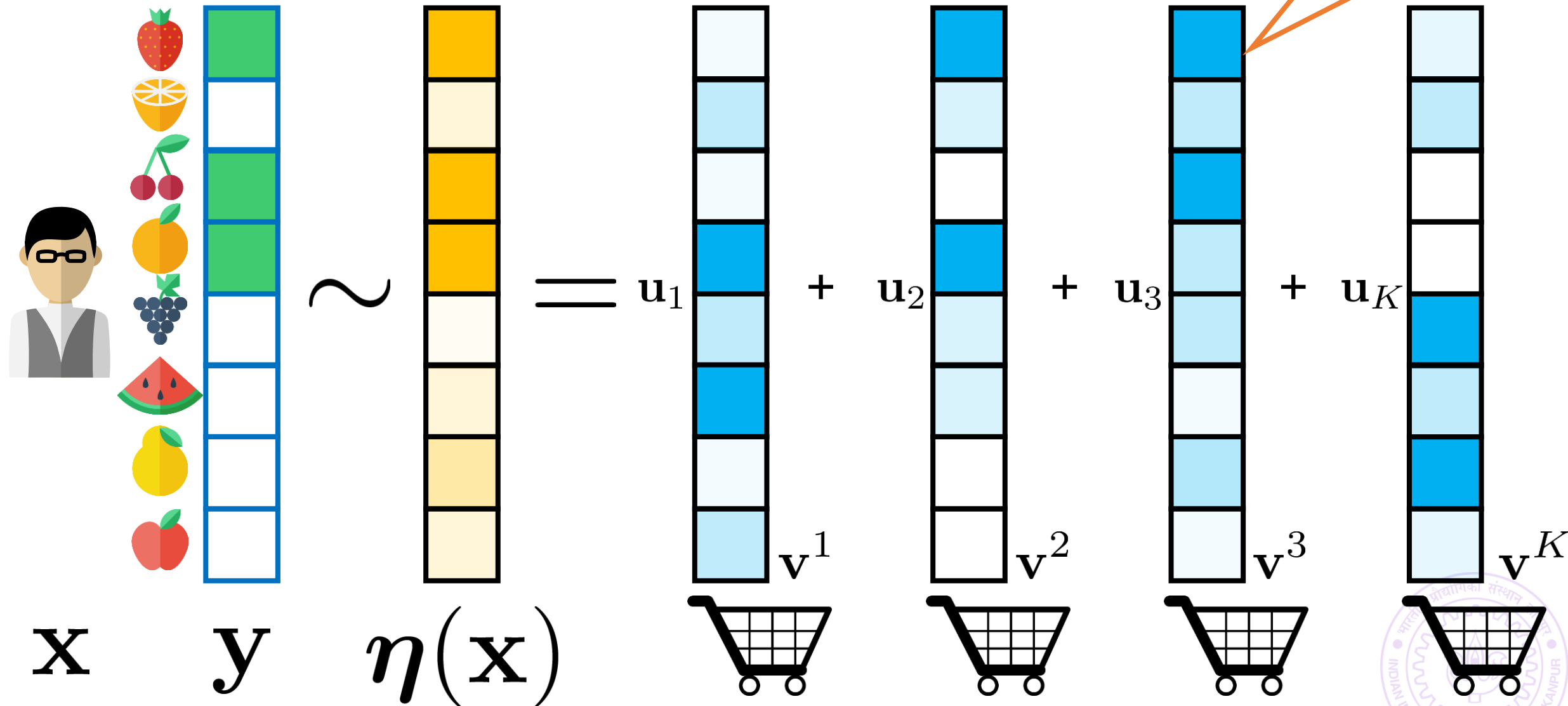
Solve as L independent Binary problems using Logistic regression!

Does not scale!

Multi-label Learning using PML



Multi-label Learning using PML



Multi-label Learning using PML

$$\mathbf{y}^i \sim \text{Bernoulli}(\boldsymbol{\eta}^i)$$

$$\boldsymbol{\eta}^i = \mathbf{V}\mathbf{u}^i, \mathbf{V} = [\mathbf{v}^1, \mathbf{v}^2, \dots, \mathbf{v}^K]$$

$$\mathbf{u}_j^i \sim \mathcal{N}(\langle \mathbf{w}^j, \mathbf{x}^i \rangle, \sigma^2), j = 1, \dots, K$$

Can fix \mathbf{V} or else ...

$$\mathbf{v}^j \sim \text{Dirichlet}(\boldsymbol{\alpha})$$

$$\mathbf{w}^j \sim \mathcal{N}(\mathbf{0}, \rho^2 \cdot I), j = 1, \dots, K$$

Generalization of Beta distribution

$$\mathbb{P}[\mathbf{v}; \boldsymbol{\alpha}] = \frac{1}{B(\boldsymbol{\alpha})} \prod_{i=1}^L \mathbf{v}_i^{\alpha_i - 1}$$

$$\mathbb{P}[p; \alpha, \beta] = \frac{1}{B(\alpha, \beta)} p^{\alpha-1} (1-p)^{\beta-1}$$

Only K models,
not L

Have fun posterioring!

Please give your Feedback

<http://tinyurl.com/ml17-18afb>