

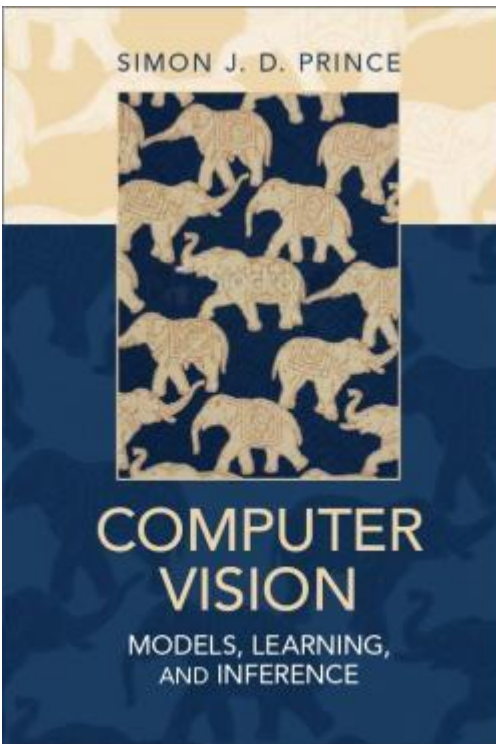
The logo of the University of Bonn, featuring a blue square in the top-left corner and a grey square in the top-right corner, separated by a white curved line.

UNIVERSITÄT **BONN**

Juergen Gall

Classification

MA-INF 2213 - Advanced Computer Vision  
SS25



## Chapter 9 Classification Models

S. Prince. **Computer Vision: Models, Learning, and Inference.** Cambridge University Press 2012

# Example application: Dogs vs. Cats



# Learn relationship

Probabilistic: Given feature  $\mathbf{x}$ , model class  $w \in \{c_1, c_2, \dots, c_k\}$  as probability distribution  $\Pr(w|\mathbf{x})$

How to model  $\Pr(\mathbf{w}|\mathbf{x})$ ?

- Choose an appropriate form for prior  $\Pr(\mathbf{w})$
- Parameterize a function of type  $\mathbf{w} = f(\mathbf{x}; \theta)$

**Learning algorithm:** learn parameters  $\theta$  from training data  $\mathbf{x}, w$

**Inference algorithm:** just evaluate  $\Pr(w|\mathbf{x})$

Similar to regression, but  $w$  has changed!

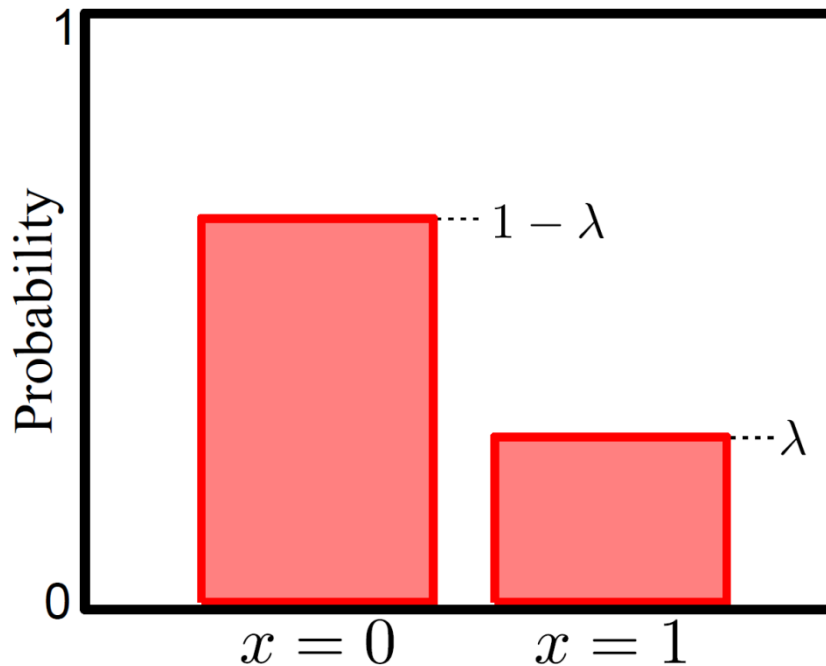
# Logistic Regression

Consider two class problem.

- Choose Bernoulli distribution over world.
- Make parameter  $\lambda$  a function of  $\mathbf{x}$

$$Pr(w|\phi_0, \phi, \mathbf{x}) = \text{Bern}_w [\text{sig}[a]]$$

# Bernoulli Distribution



$$Pr(x = 0) = 1 - \lambda$$

$$Pr(x = 1) = \lambda.$$

or

$$Pr(x) = \lambda^x (1 - \lambda)^{1-x}$$

For short we write:

$$Pr(x) = \text{Bern}_x[\lambda]$$

Bernoulli distribution describes situation where only two possible outcomes  $x=0/x=1$  or failure/success

Takes a single parameter  $\lambda \in [0, 1]$

# Logistic Regression

Consider two class problem.

- Choose Bernoulli distribution over world.
- Make parameter  $\lambda$  a function of  $\mathbf{x}$

$$Pr(w|\phi_0, \phi, \mathbf{x}) = \text{Bern}_w [\text{sig}[a]]$$

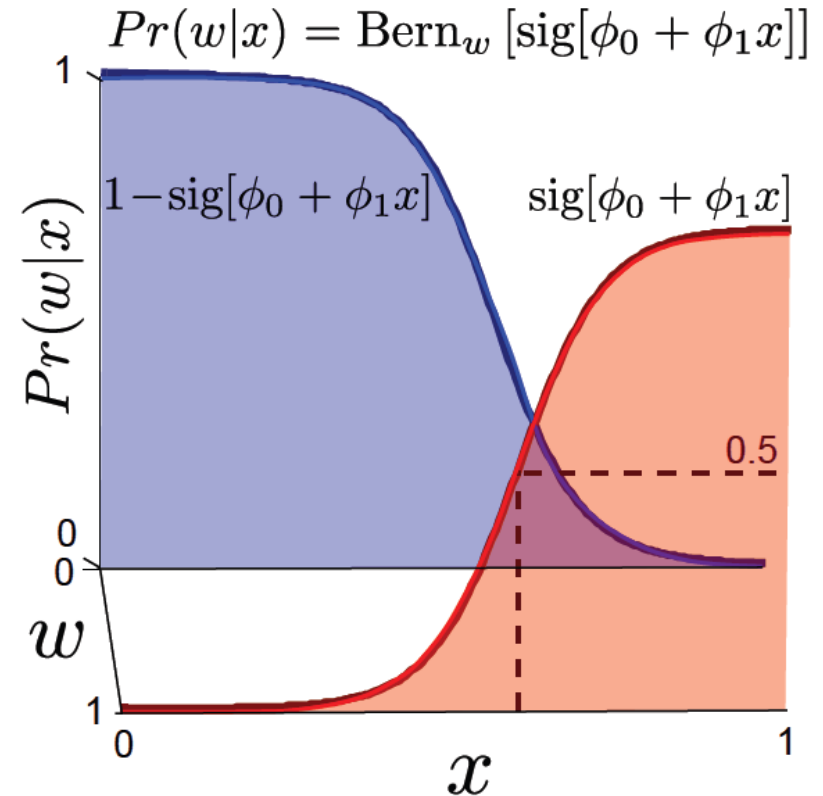
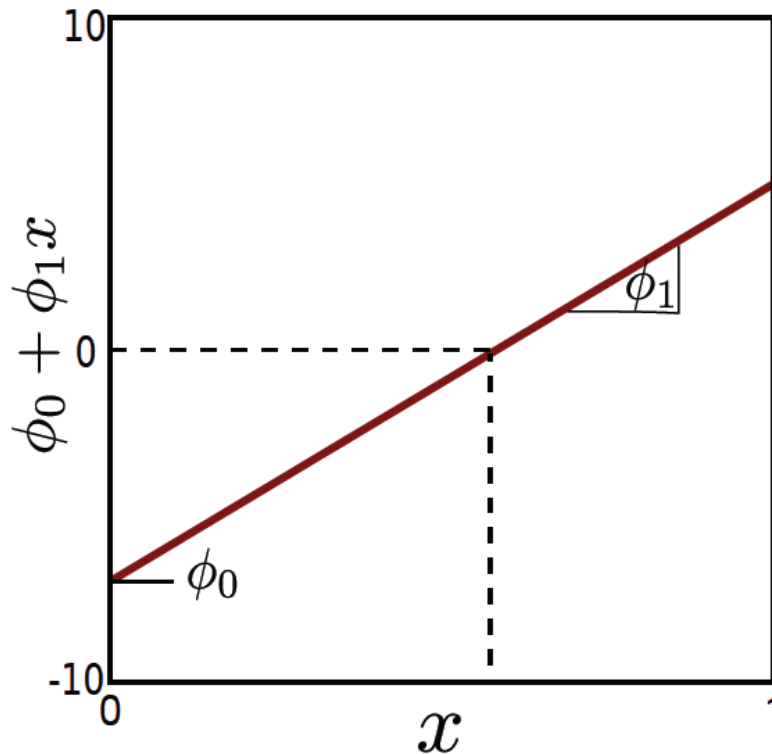
Model **activation** with a linear function

$$a = \phi_0 + \phi^T \mathbf{x}$$

creates number between  $[-\infty, \infty]$ . Maps to  $[0, 1]$  with

$$\text{sig}[a] = \frac{1}{1 + \exp[-a]}$$

# Logistic Regression



Two parameters  $\theta = \{\phi_0, \phi_1\}$

Learning by standard methods (ML, MAP, Bayesian)

Inference: Just evaluate  $Pr(w|x)$



# Neater Notation

$$Pr(w|\phi_0, \phi, \mathbf{x}) = \text{Bern}_w [\text{sig}[a]]$$

To make notation easier to handle, we

- Attach a 1 to the start of every data vector

$$\mathbf{x}_i \leftarrow [1 \quad \mathbf{x}_i^T]^T$$

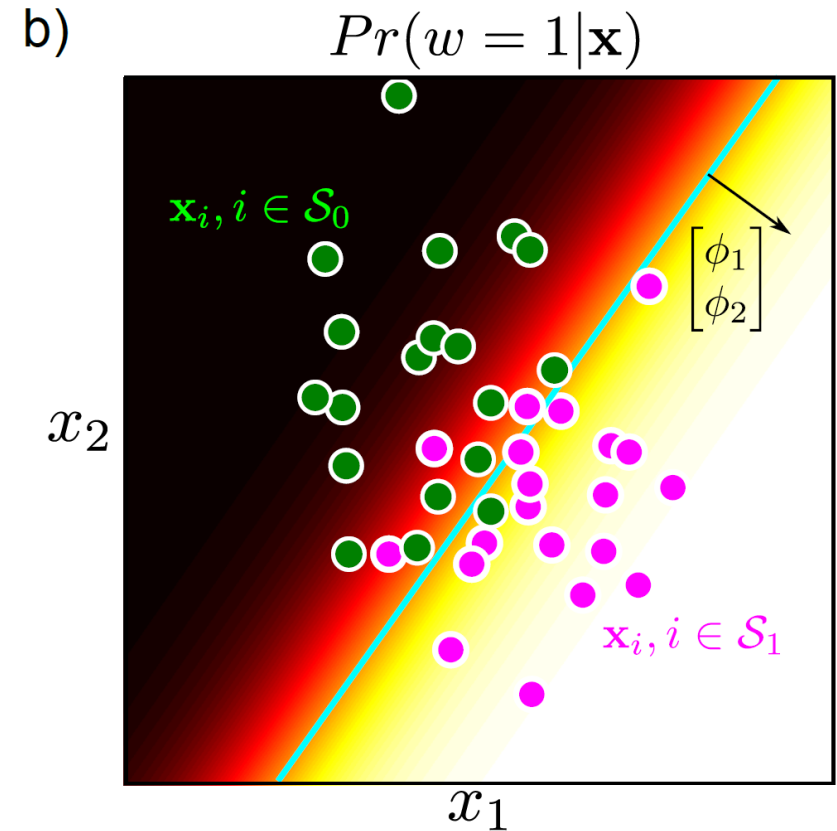
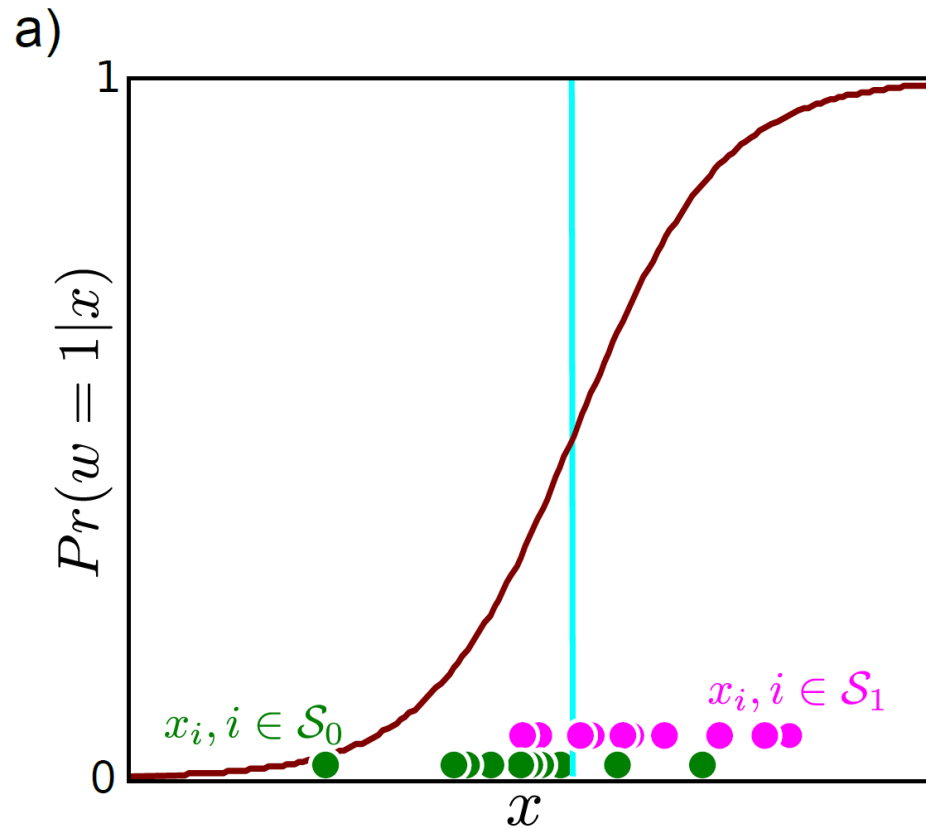
- Attach the offset to the start of the gradient vector  $\phi$

$$\phi \leftarrow [\phi_0 \quad \phi^T]^T$$

New model:

$$Pr(w|\phi, \mathbf{x}) = \text{Bern}_w \left[ \frac{1}{1 + \exp[-\phi^T \mathbf{x}]} \right]$$

# Logistic regression



$$Pr(w|\phi, \mathbf{x}) = \text{Bern}_w \left[ \frac{1}{1 + \exp[-\phi^T \mathbf{x}]} \right]$$

# Maximum Likelihood

$$\begin{aligned} Pr(\mathbf{w}|\mathbf{X}, \phi) &= \prod_{i=1}^I \lambda^{w_i} (1 - \lambda)^{1-w_i} \\ &= \prod_{i=1}^I \left( \frac{1}{1 + \exp[-\phi^T \mathbf{x}_i]} \right)^{w_i} \left( \frac{\exp[-\phi^T \mathbf{x}_i]}{1 + \exp[-\phi^T \mathbf{x}_i]} \right)^{1-w_i} \end{aligned}$$

Take logarithm

$$L = \sum_{i=1}^I w_i \log \left[ \frac{1}{1 + \exp[-\phi^T \mathbf{x}_i]} \right] + \sum_{i=1}^I (1 - w_i) \log \left[ \frac{\exp[-\phi^T \mathbf{x}_i]}{1 + \exp[-\phi^T \mathbf{x}_i]} \right]$$

Take derivative:

$$\frac{\partial L}{\partial \phi} = - \sum_{i=1}^I \left( \frac{1}{1 + \exp[-\phi^T \mathbf{x}_i]} - w_i \right) \mathbf{x}_i = - \sum_{i=1}^I (\text{sig}[a_i] - w_i) \mathbf{x}_i$$

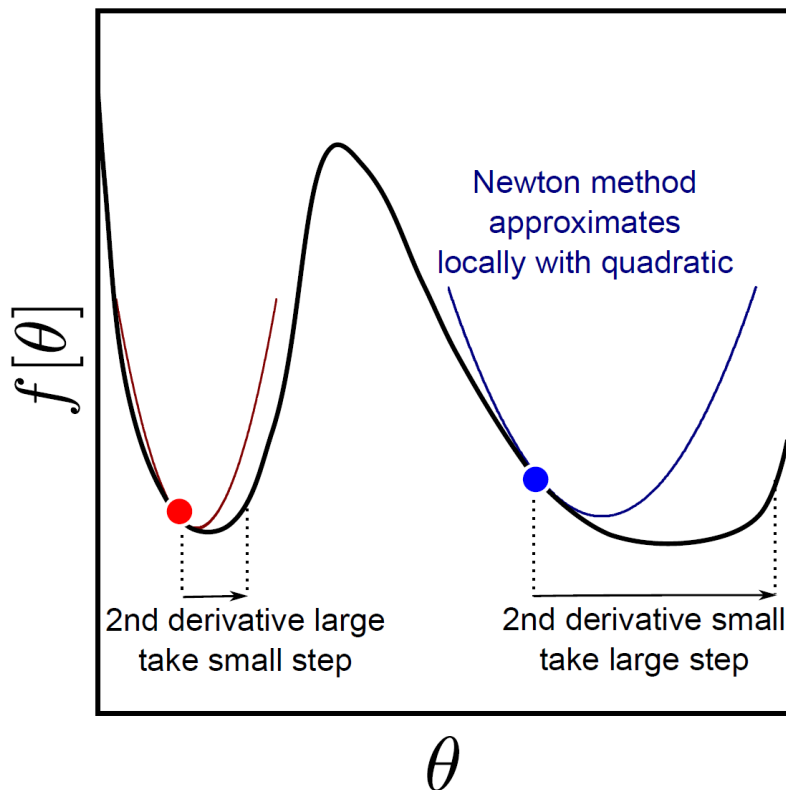
$$\frac{\partial L}{\partial \phi} = - \sum_{i=1}^I \left( \frac{1}{1 + \exp[-\phi^T \mathbf{x}_i]} - w_i \right) \mathbf{x}_i = - \sum_{i=1}^I (\text{sig}[a_i] - w_i) \mathbf{x}_i$$

Unfortunately, there is no closed form solution– we cannot get an expression for  $\phi$  in terms of  $\mathbf{x}$  and  $w$

Use iterative non-linear optimization

# Newton's Method

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} [f[\theta]] \quad \theta^{[t+1]} = \theta^{[t]} - \lambda \left( \frac{\partial^2 f}{\partial \theta^2} \right)^{-1} \frac{\partial f}{\partial \theta}$$



Matrix of second derivatives is called the Hessian.

If positive definite, then convex

# Optimization for Logistic Regression

$$\phi^{[t]} = \phi^{[t-1]} + \alpha \left( \frac{\partial^2 L}{\partial \phi^2} \right)^{-1} \frac{\partial L}{\partial \phi}$$

Derivatives of log likelihood:

$$\frac{\partial L}{\partial \phi} = - \sum_{i=1}^I (\text{sig}[a_i] - w_i) \mathbf{x}_i$$

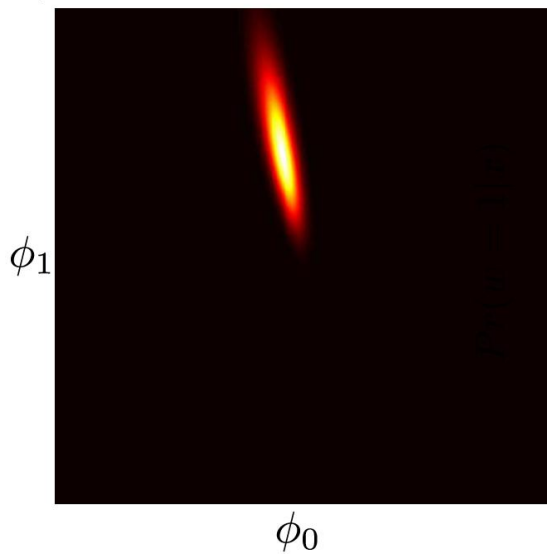
$$\frac{\partial^2 L}{\partial \phi^2} = - \sum_{i=1}^I \text{sig}[a_i] (1 - \text{sig}[a_i]) \mathbf{x}_i \mathbf{x}_i^T$$



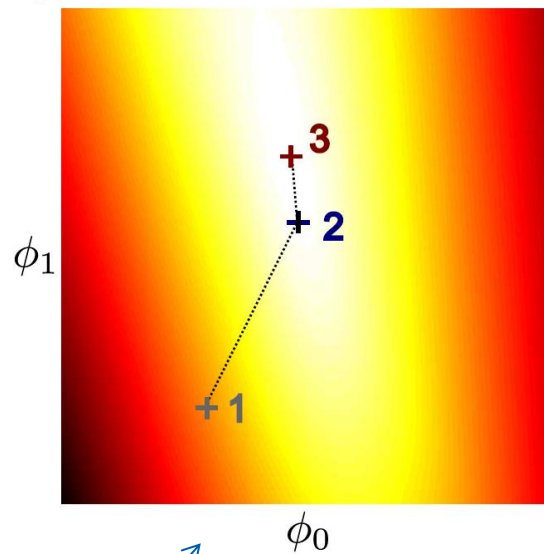
Positive definite!

$$Pr(\mathbf{w}|\mathbf{X}, \phi) = \prod_{i=1}^I \left( \frac{1}{1 + \exp[-\phi^T \mathbf{x}_i]} \right)^{w_i} \left( \frac{\exp[-\phi^T \mathbf{x}_i]}{1 + \exp[-\phi^T \mathbf{x}_i]} \right)^{1-w_i}$$

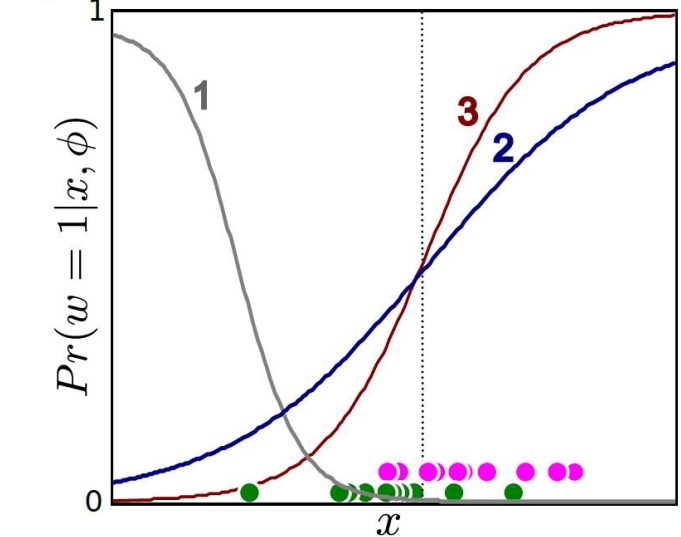
a)



b)

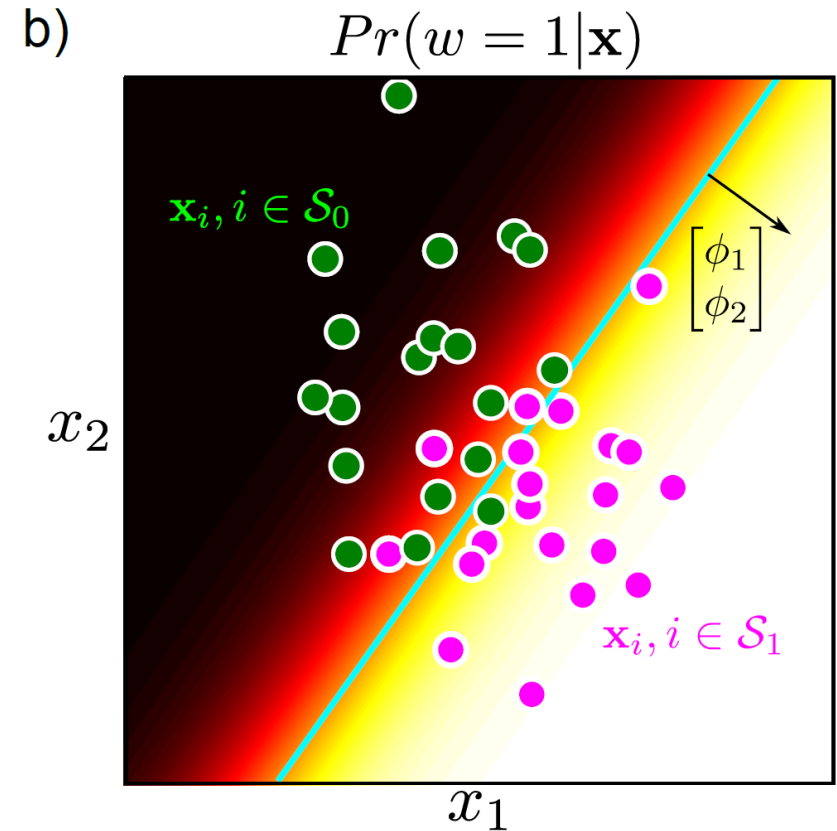
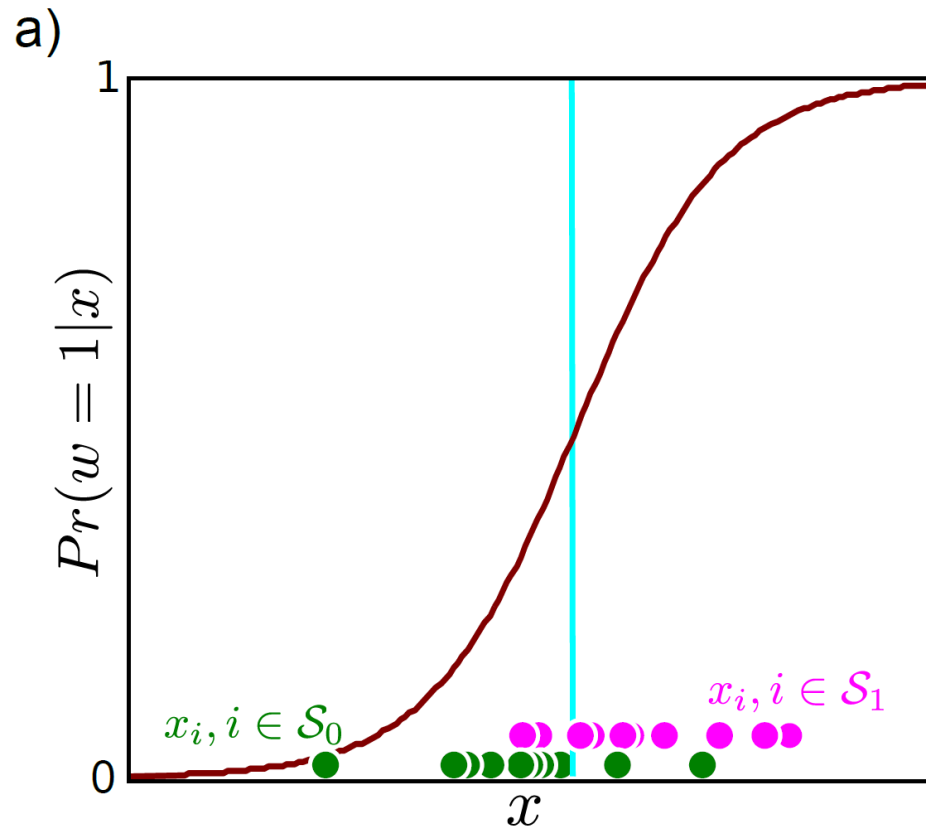


c)



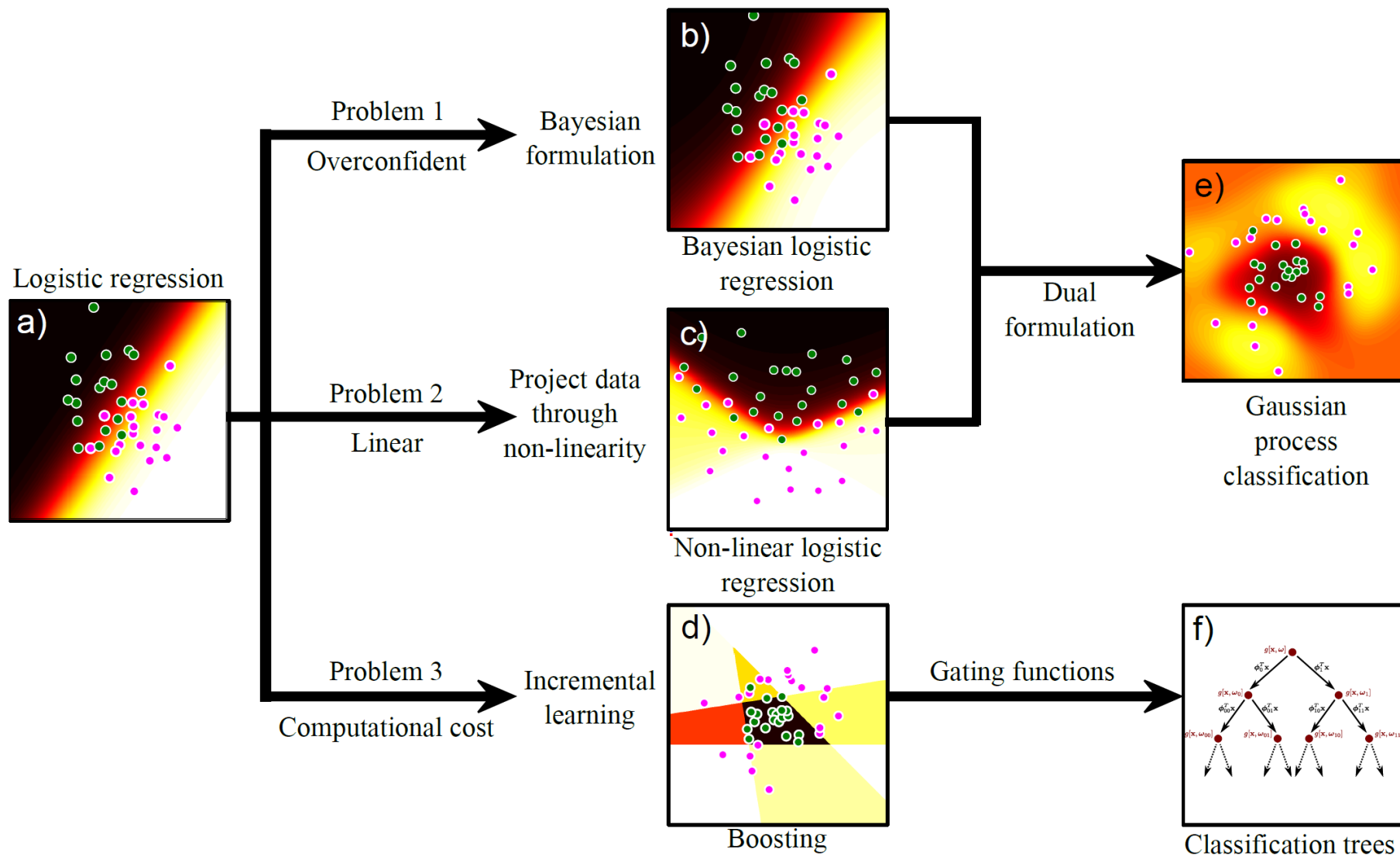
$$L = \sum_{i=1}^I w_i \log \left[ \frac{1}{1 + \exp[-\phi^T \mathbf{x}_i]} \right] + \sum_{i=1}^I (1 - w_i) \log \left[ \frac{\exp[-\phi^T \mathbf{x}_i]}{1 + \exp[-\phi^T \mathbf{x}_i]} \right]$$

# Maximum likelihood fits



$$Pr(w|\phi, \mathbf{x}) = \text{Bern}_w \left[ \frac{1}{1 + \exp[-\phi^T \mathbf{x}]} \right]$$





# Bayesian Logistic Regression

Likelihood:

$$Pr(\mathbf{w}|\mathbf{X}, \phi) = \prod_{i=1}^I \left( \frac{1}{1 + \exp[-\phi^T \mathbf{x}_i]} \right)^{w_i} \left( \frac{\exp[-\phi^T \mathbf{x}_i]}{1 + \exp[-\phi^T \mathbf{x}_i]} \right)^{1-w_i}$$

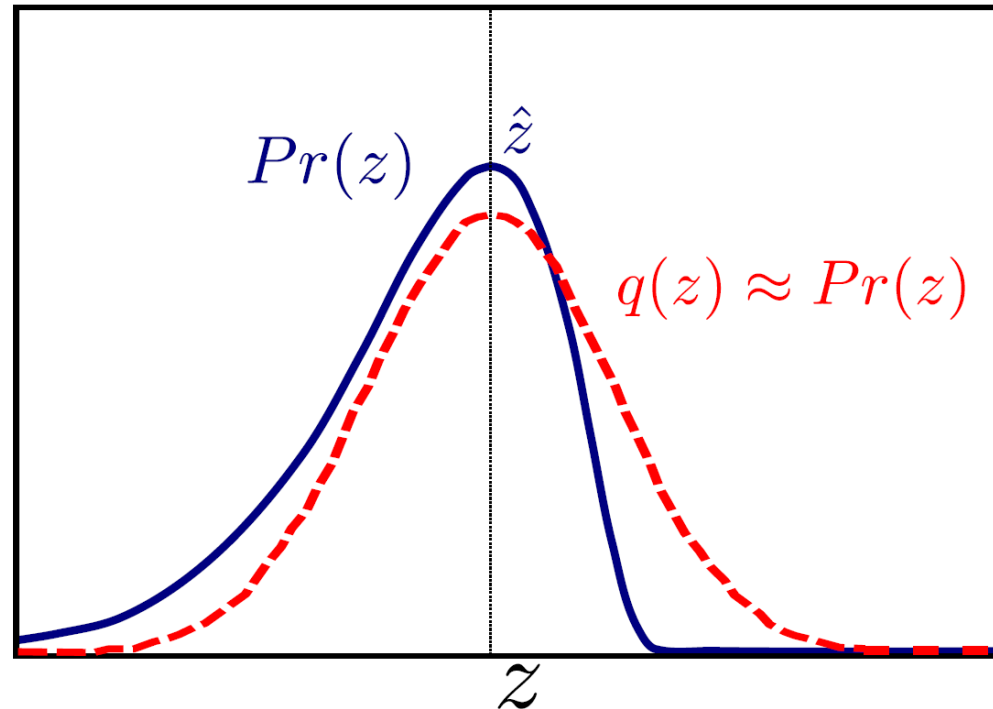
Prior:

$$Pr(\phi) = \text{Norm}_{\phi}[\mathbf{0}, \sigma_p^2 \mathbf{I}]$$

Apply Bayes' rule:

$$Pr(\phi|\mathbf{X}, \mathbf{w}) = \frac{Pr(\mathbf{w}|\mathbf{X}, \phi) Pr(\phi)}{Pr(\mathbf{w}|\mathbf{X})}$$

# Laplace Approximation



Approximate posterior distribution with normal

- Set mean to MAP estimate
- Set covariance to match that at MAP estimate (actually: get 2<sup>nd</sup> derivatives to agree)

# Laplace Approximation

Find MAP solution by optimizing

$$L = \sum_{i=1}^I \log[Pr(w_i | \mathbf{x}_i, \phi)] + \log[Pr(\phi)]$$

using Newton's method:

$$\begin{aligned} \frac{\partial L}{\partial \phi} &= - \sum_{i=1}^I (\text{sig}[a_i] - w_i) \mathbf{x}_i - \frac{\phi}{\sigma_p^2} \\ \frac{\partial^2 L}{\partial \phi^2} &= - \sum_{i=1}^I \text{sig}[a_i] (1 - \text{sig}[a_i]) \mathbf{x}_i \mathbf{x}_i^T - \frac{1}{\sigma_p^2} \end{aligned}$$

# Laplace Approximation

Find MAP solution by optimizing

$$L = \sum_{i=1}^I \log[Pr(w_i | \mathbf{x}_i, \phi)] + \log[Pr(\phi)]$$

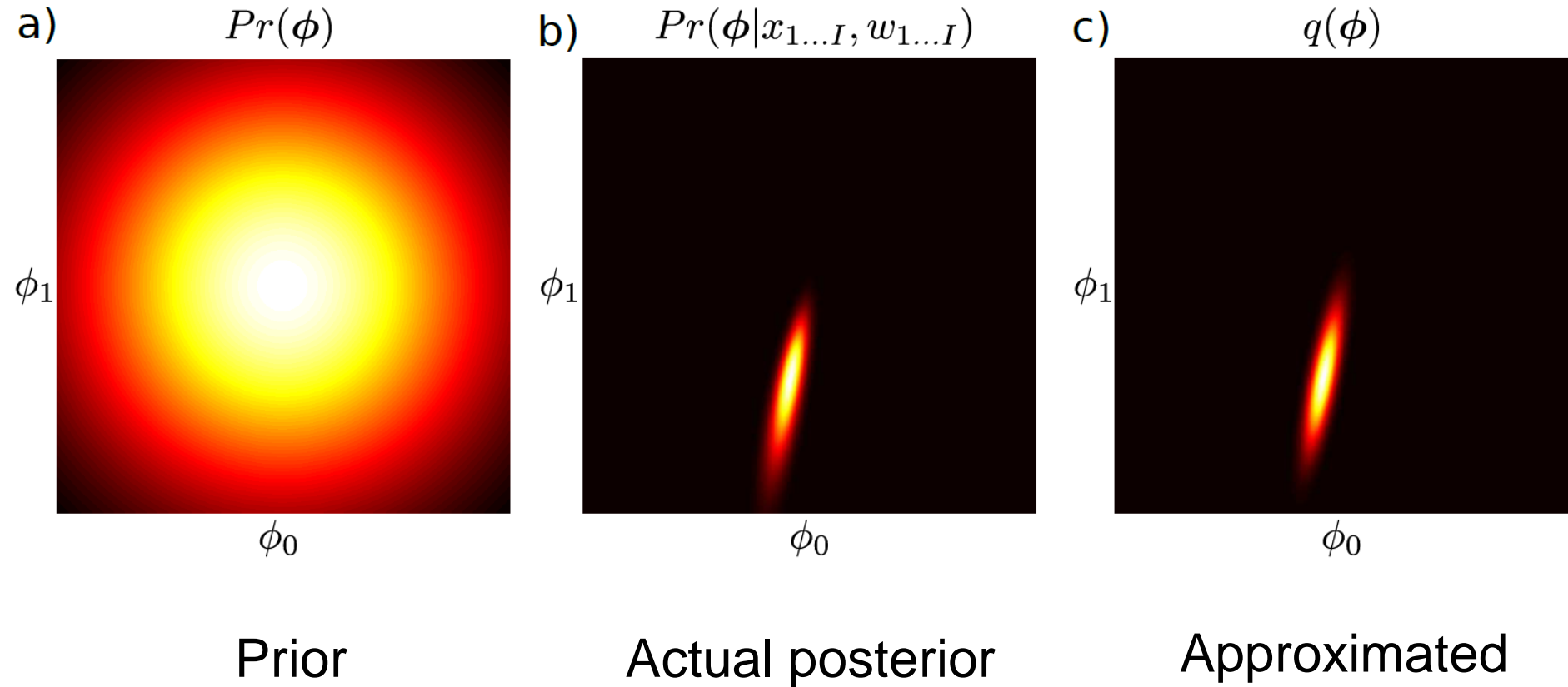
Approximate with normal

$$Pr(\phi | \mathbf{X}, \mathbf{w}) \approx q(\phi) = \text{Norm}_{\phi}[\boldsymbol{\mu}, \boldsymbol{\Sigma}]$$

where

$$\begin{aligned} \boldsymbol{\mu} &= \hat{\phi} \\ \boldsymbol{\Sigma} &= - \left( \frac{\partial^2 L}{\partial \phi^2} \right)^{-1} \bigg|_{\phi = \hat{\phi}} \end{aligned}$$

# Laplace Approximation



$$\begin{aligned} Pr(w^* | \mathbf{x}^*, \mathbf{X}, \mathbf{w}) &= \int Pr(w^* | \mathbf{x}^*, \phi) Pr(\phi | \mathbf{X}, \mathbf{w}) d\phi \\ &\approx \int Pr(w^* | \mathbf{x}^*, \phi) q(\phi) d\phi. \end{aligned}$$

Can re-express in terms of activation  $a = \phi^T \mathbf{x}^*$

$$Pr(w^* | \mathbf{x}^*, \mathbf{X}, \mathbf{w}) \approx \int Pr(w^* | a) Pr(a) da$$

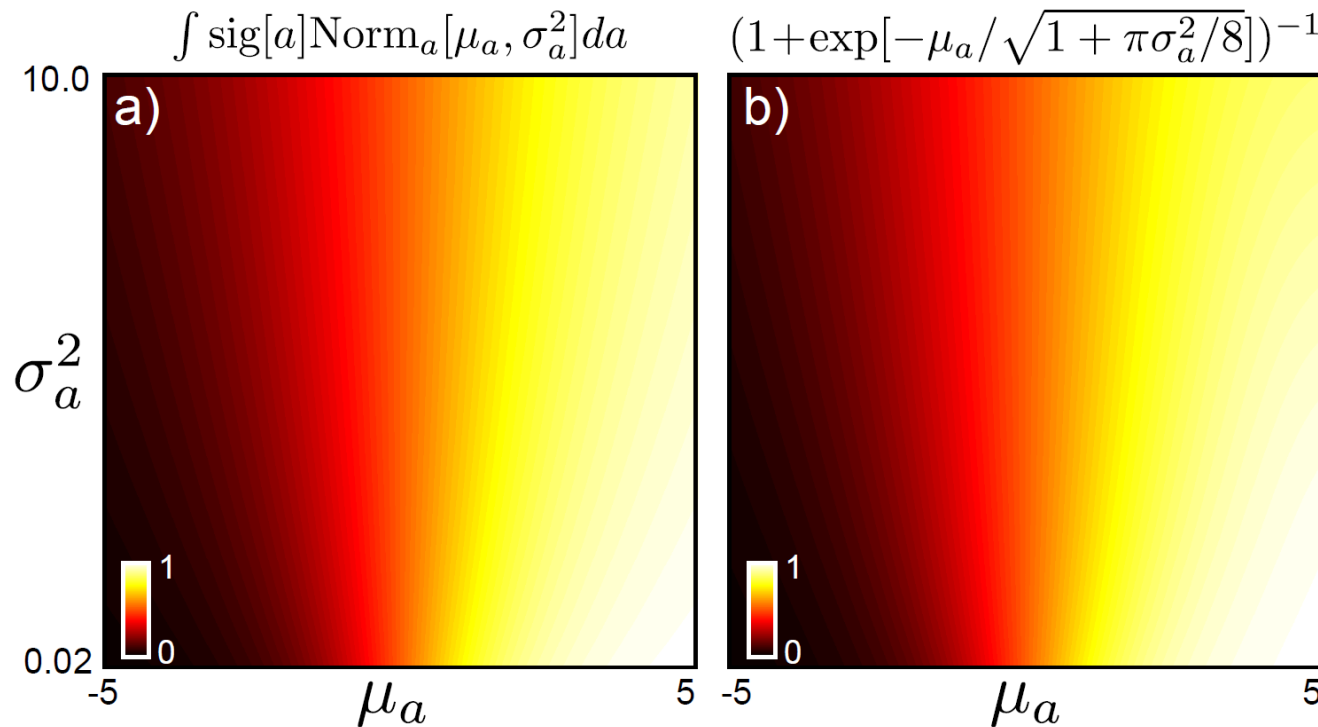
Using transformation properties of normal distributions

$$\begin{aligned} Pr(a) = Pr(\phi^T \mathbf{x}^*) &= \text{Norm}_a[\boldsymbol{\mu}^T \mathbf{x}^*, \mathbf{x}^{*T} \boldsymbol{\Sigma} \mathbf{x}^*] \\ &= \text{Norm}_a[\mu_a, \sigma_a^2], \end{aligned}$$

# Approximation of Integral

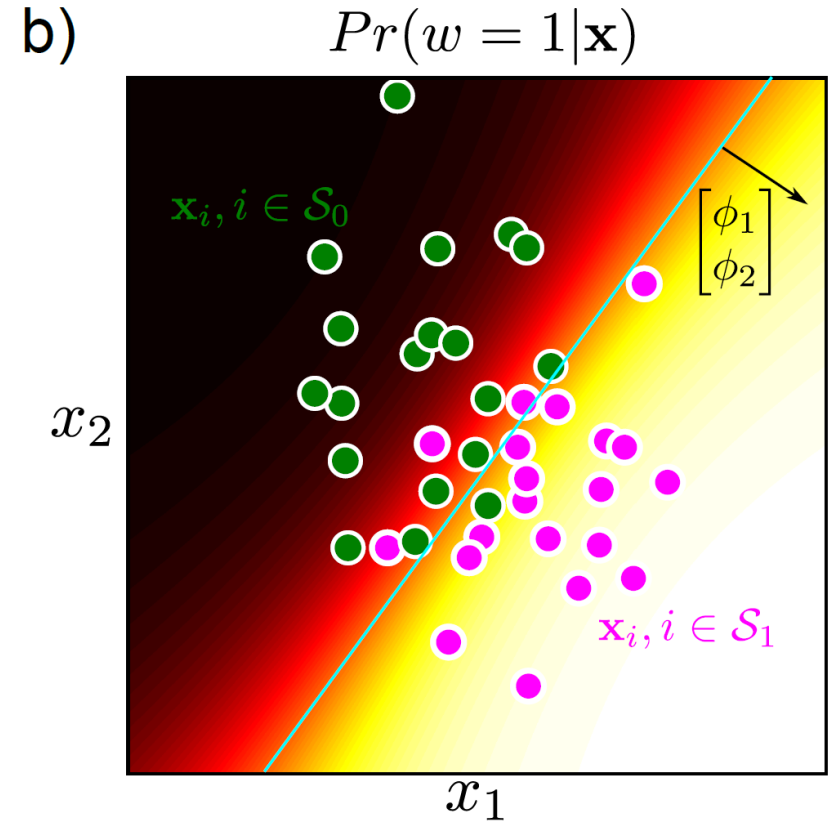
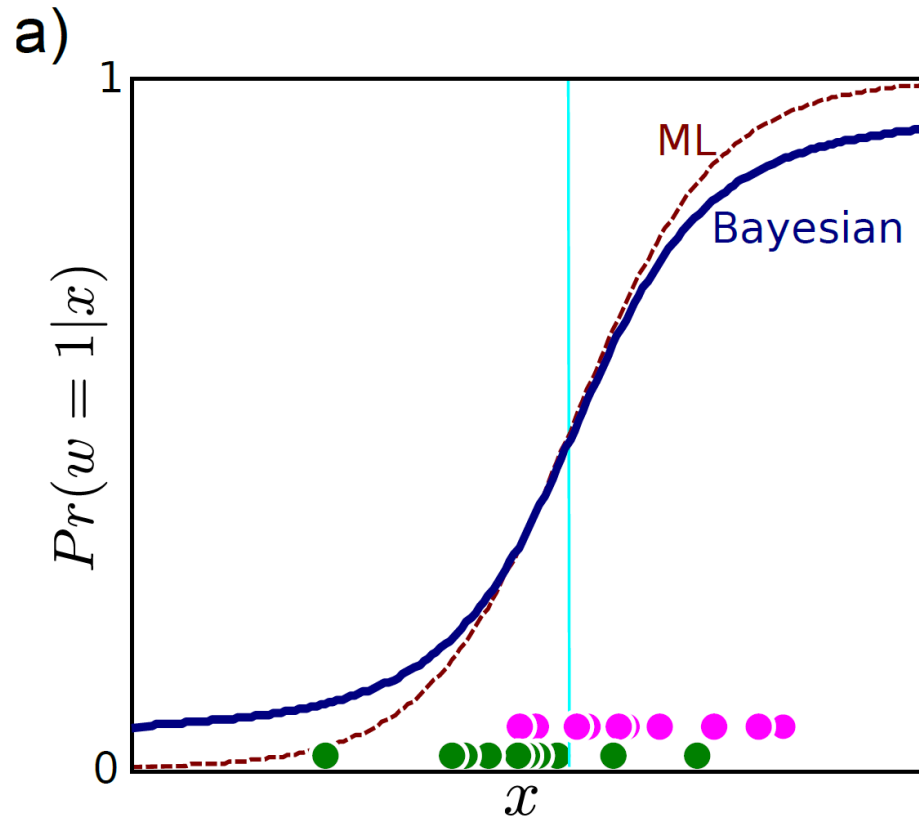
(Or perform numerical integration on  $a$  – which is 1D)

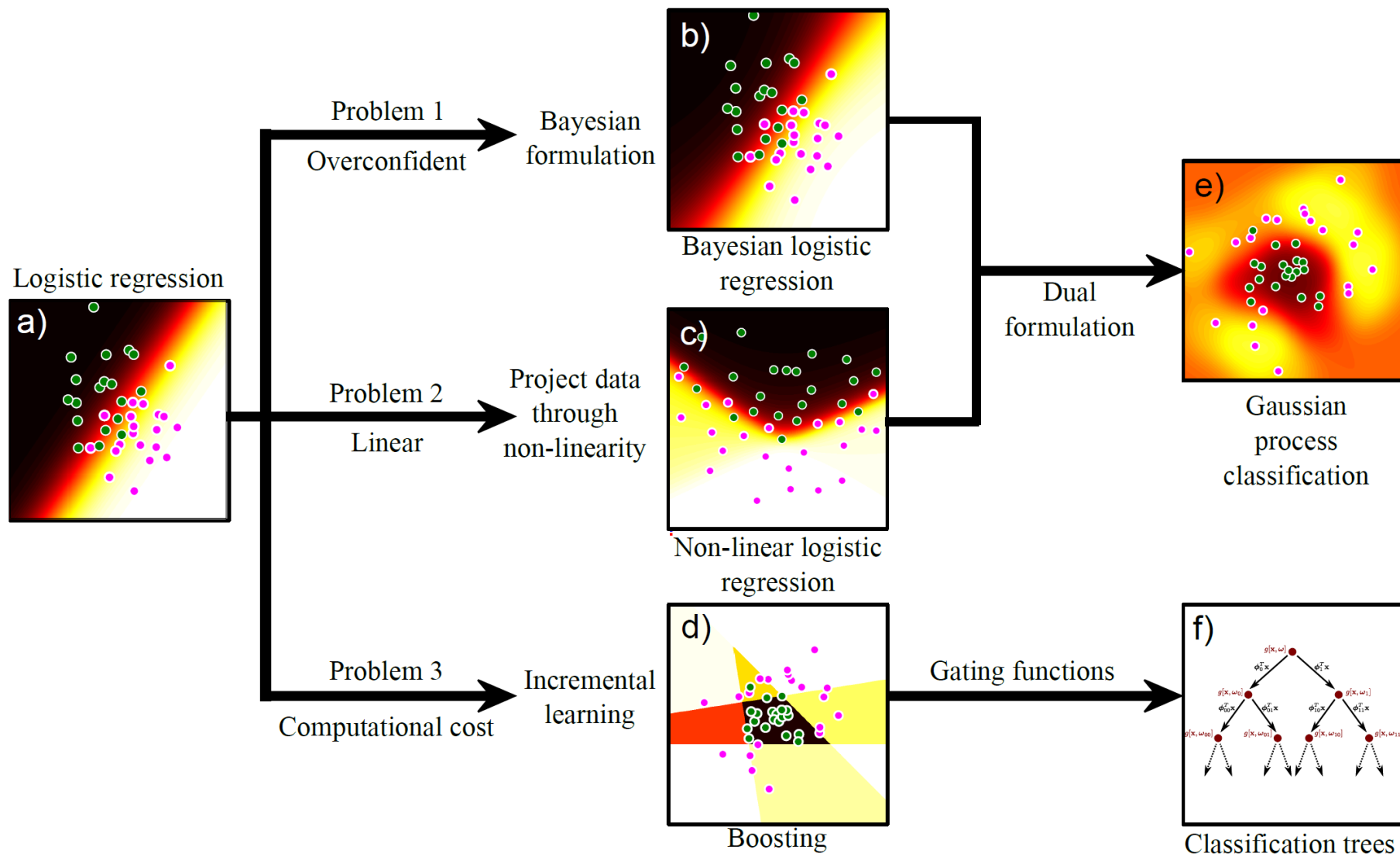
$$\int Pr(w^*|a) \text{Norm}_a[\mu_a, \sigma_a^2] da \approx \frac{1}{1 + \exp[-\mu_a / \sqrt{1 + \pi \sigma_a^2 / 8}]}$$





# Bayesian Solution





# Non-linear logistic regression

Same idea as for regression.

- Apply non-linear transformation

$$\mathbf{z} = \mathbf{f}[\mathbf{x}]$$

- Build model as usual

$$\begin{aligned} Pr(w = 1 | \mathbf{x}, \phi) &= \text{Bern}_w \left[ \text{sig}[\phi^T \mathbf{z}] \right] \\ &= \text{Bern}_w \left[ \text{sig}[\phi^T \mathbf{f}[\mathbf{x}]] \right] \end{aligned}$$

# Non-linear logistic regression

## Example transformations:

- Arc tan functions of projections:  $z_k = \arctan[\boldsymbol{\alpha}_k^T \mathbf{x}]$
- Radial basis functions:  $z_k = \exp \left[ -\frac{1}{\lambda_0} (\mathbf{x} - \boldsymbol{\alpha}_k)^T (\mathbf{x} - \boldsymbol{\alpha}_k) \right]$

# Non-linear logistic regression

## Example transformations:

- Arc tan functions of projections:  $z_k = \arctan[\boldsymbol{\alpha}_k^T \mathbf{x}]$
- Radial basis functions:  $z_k = \exp \left[ -\frac{1}{\lambda_0} (\mathbf{x} - \boldsymbol{\alpha}_k)^T (\mathbf{x} - \boldsymbol{\alpha}_k) \right]$

## Fit using optimization (also transformation parameters $\boldsymbol{\alpha}$ ):

$$\boldsymbol{\theta} = [\boldsymbol{\phi}^T, \boldsymbol{\alpha}_1^T, \boldsymbol{\alpha}_2^T, \dots, \boldsymbol{\alpha}_K^T]^T \quad a_i = \boldsymbol{\phi}^T \mathbf{f}[\mathbf{x}_i]$$

$$\frac{\partial L}{\partial \boldsymbol{\theta}} = \sum_{i=1}^I (w_i - \text{sig}[a_i]) \frac{\partial a_i}{\partial \boldsymbol{\theta}}$$

$$\frac{\partial^2 L}{\partial \boldsymbol{\theta}^2} = \sum_{i=1}^I \text{sig}[a_i] (\text{sig}[a_i] - 1) \frac{\partial a_i}{\partial \boldsymbol{\theta}} \frac{\partial a_i}{\partial \boldsymbol{\theta}}^T + (w_i - \text{sig}[a_i]) \frac{\partial^2 a_i}{\partial \boldsymbol{\theta}^2}$$

# Linear logistic regression (recall)

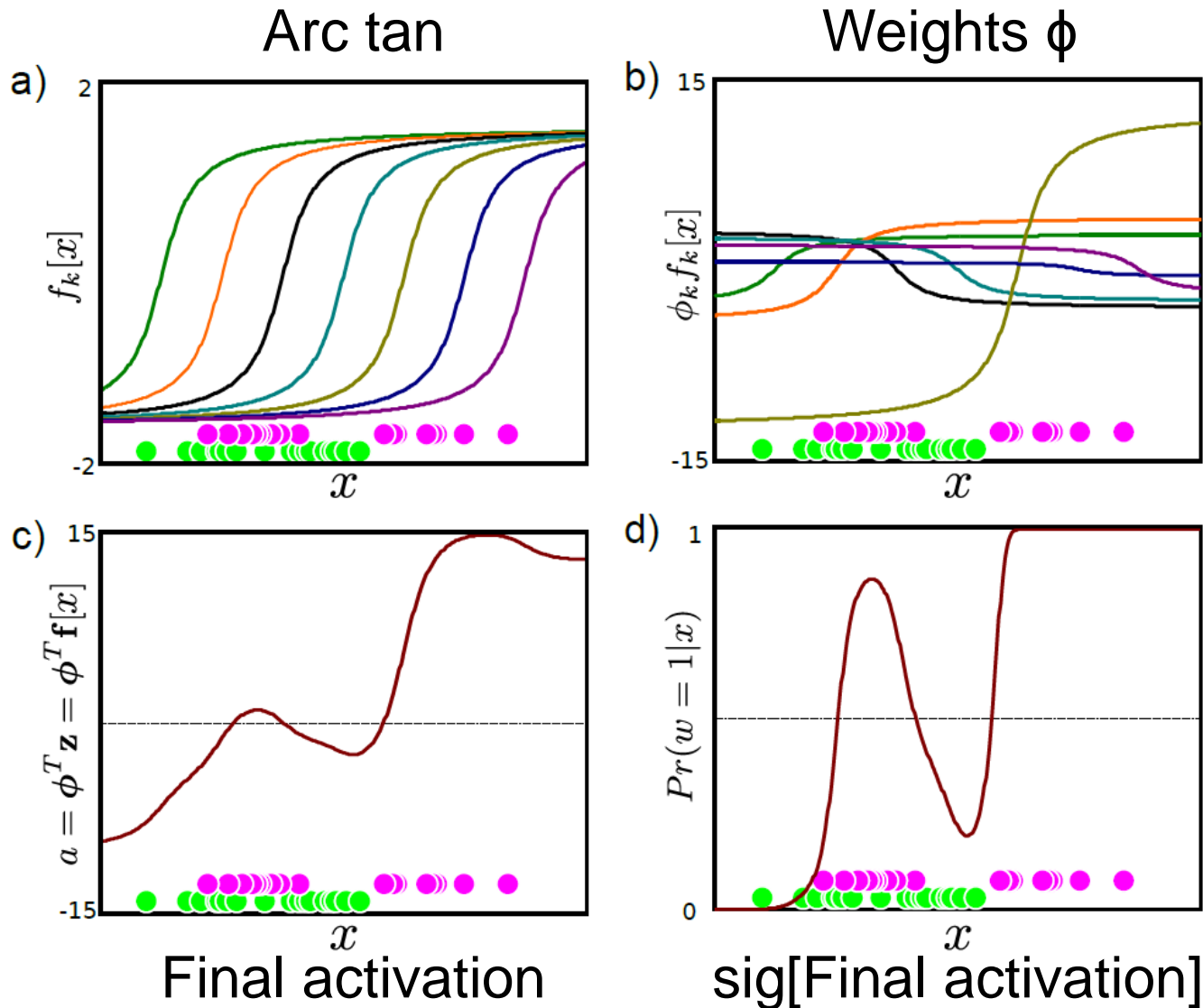
$$\phi^{[t]} = \phi^{[t-1]} + \alpha \left( \frac{\partial^2 L}{\partial \phi^2} \right)^{-1} \frac{\partial L}{\partial \phi}$$

Derivatives of log likelihood:

$$\frac{\partial L}{\partial \phi} = - \sum_{i=1}^I (\text{sig}[a_i] - w_i) \mathbf{x}_i$$

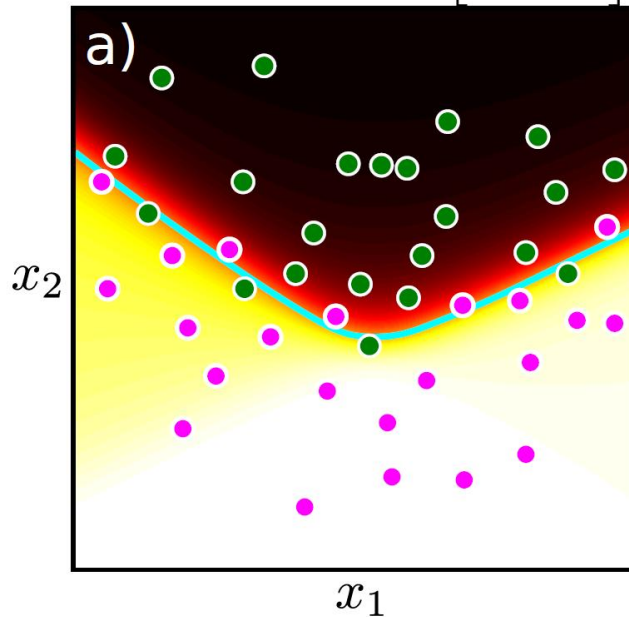
$$\frac{\partial^2 L}{\partial \phi^2} = - \sum_{i=1}^I \text{sig}[a_i] (1 - \text{sig}[a_i]) \mathbf{x}_i \mathbf{x}_i^T$$

# Non-linear logistic regression in 1D

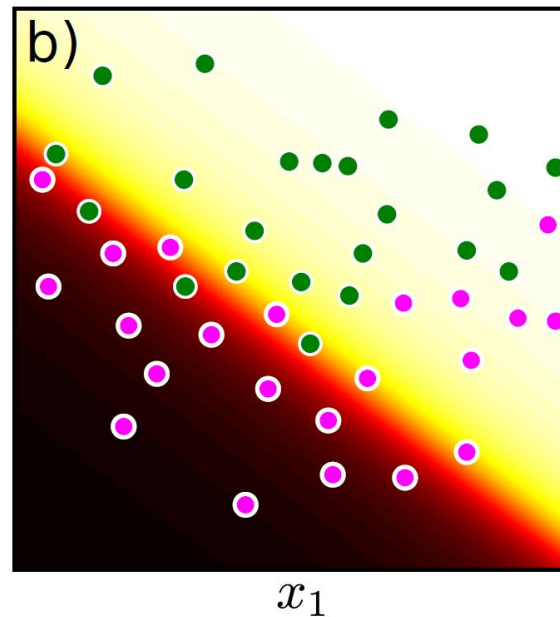


# Non-linear logistic regression in 2D

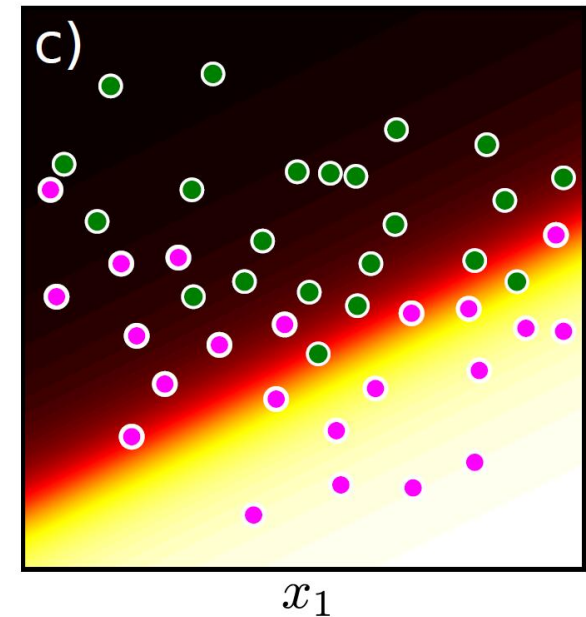
$$Pr(w=1|\mathbf{x}) = \text{sig} \left[ \phi^T \mathbf{f}[\mathbf{x}] \right]$$



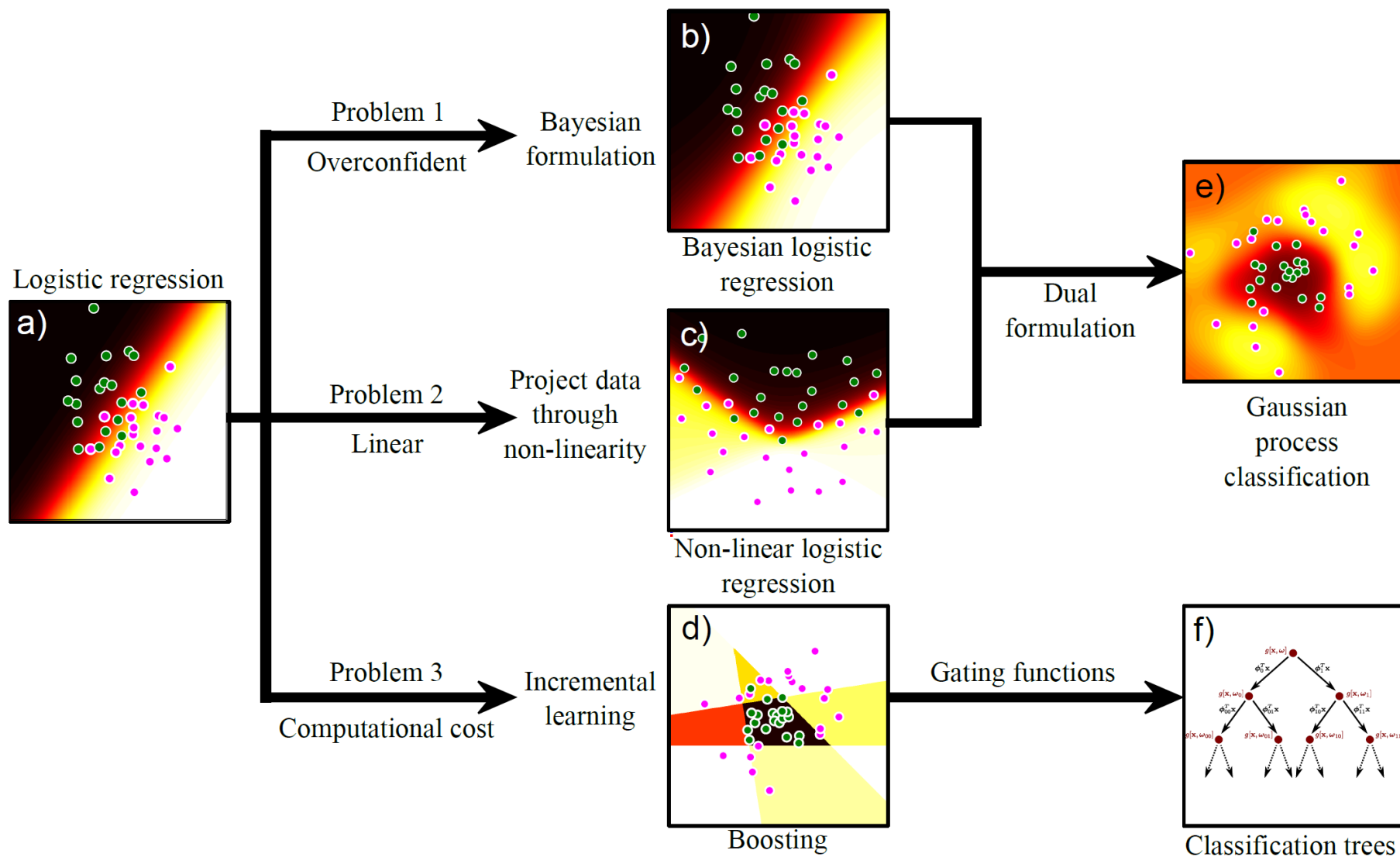
$$f_1[\mathbf{x}] = \arctan [\alpha_1^T \mathbf{x}]$$



$$f_2[\mathbf{x}] = \arctan [\alpha_2^T \mathbf{x}]$$







# Dual Logistic Regression

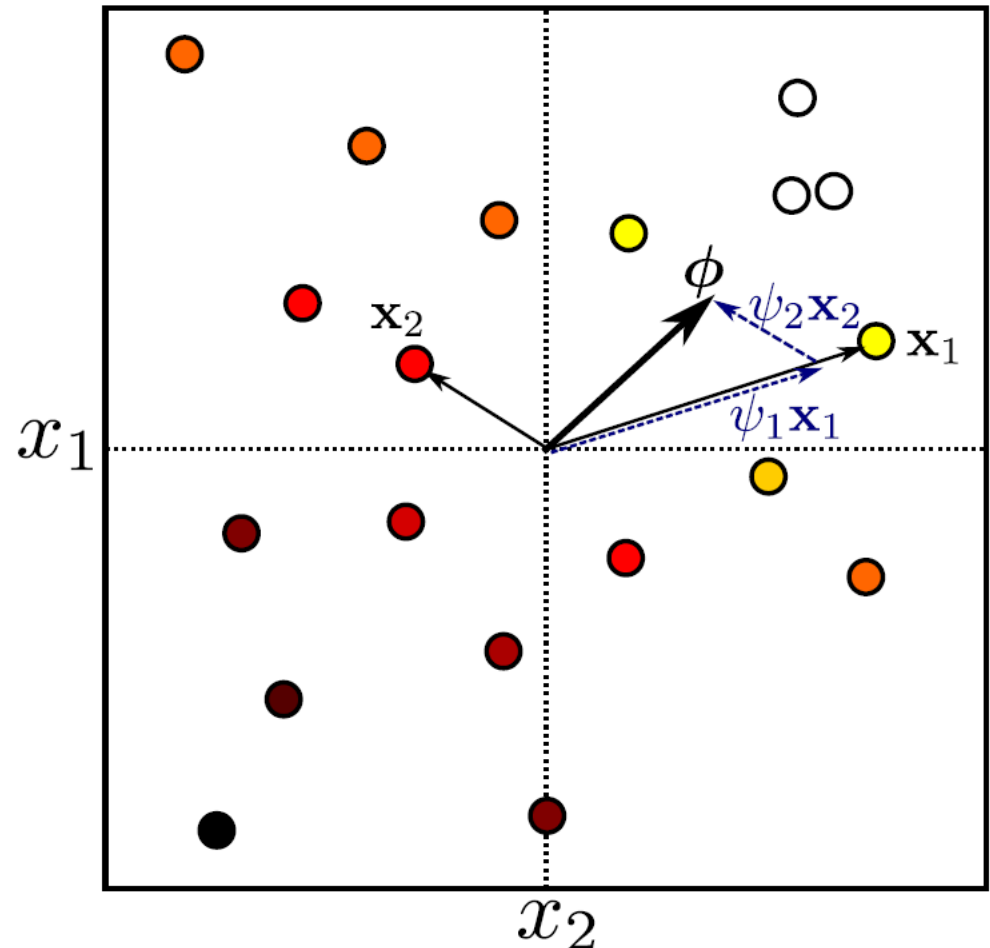
## KEY IDEA:

Gradient  $\Phi$  is just a vector in the data space

Can represent as a weighted sum of the data points

$$\phi = \mathbf{X}\psi$$

Now solve for  $\Psi$ . One parameter per training example.



# Maximum Likelihood

## Likelihood

$$Pr(\mathbf{w}|\mathbf{X}, \boldsymbol{\psi}) = \prod_{i=1}^I \text{Bern}_{w_i} [\text{sig}[a_i]] = \prod_{i=1}^I \text{Bern}_{w_i} [\text{sig}[\boldsymbol{\psi}^T \mathbf{X}^T \mathbf{x}_i]]$$

## Derivatives

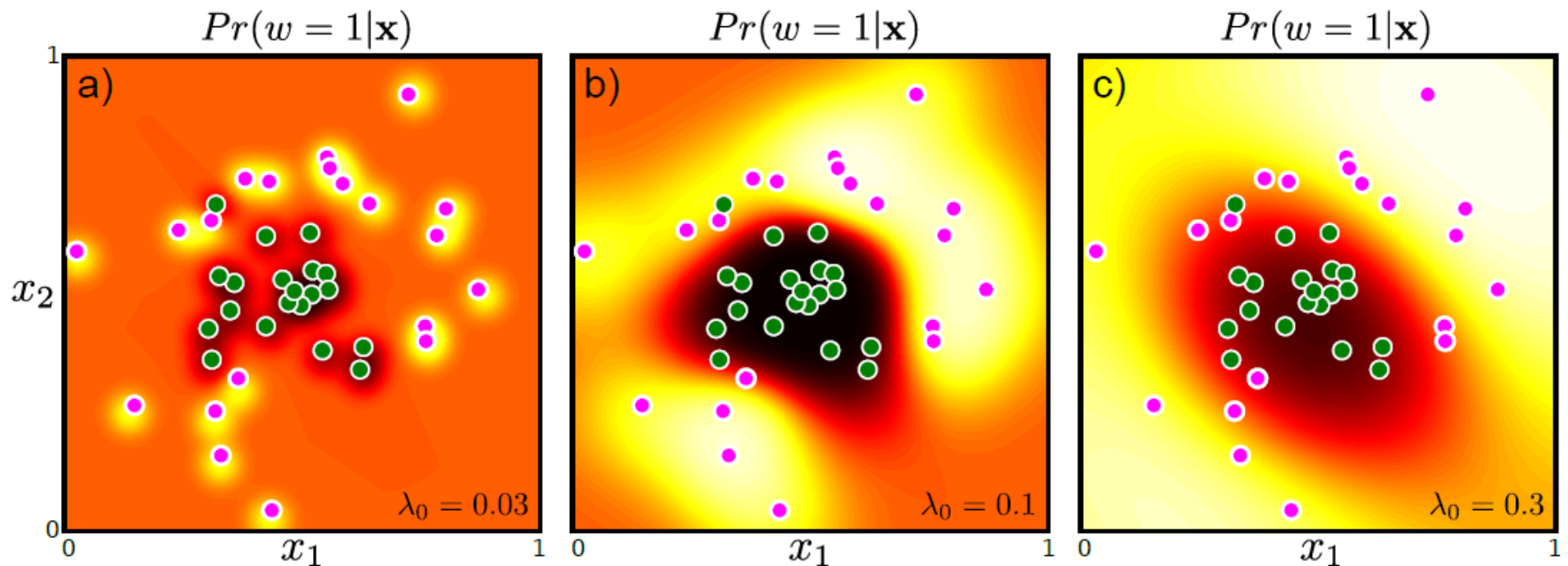
$$\frac{\partial L}{\partial \boldsymbol{\psi}} = - \sum_{i=1}^I (\text{sig}[a_i] - w_i) \mathbf{X}^T \mathbf{x}_i$$

$$\frac{\partial^2 L}{\partial \boldsymbol{\psi}^2} = - \sum_{i=1}^I \text{sig}[a_i] (1 - \text{sig}[a_i]) \mathbf{X}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{X}$$

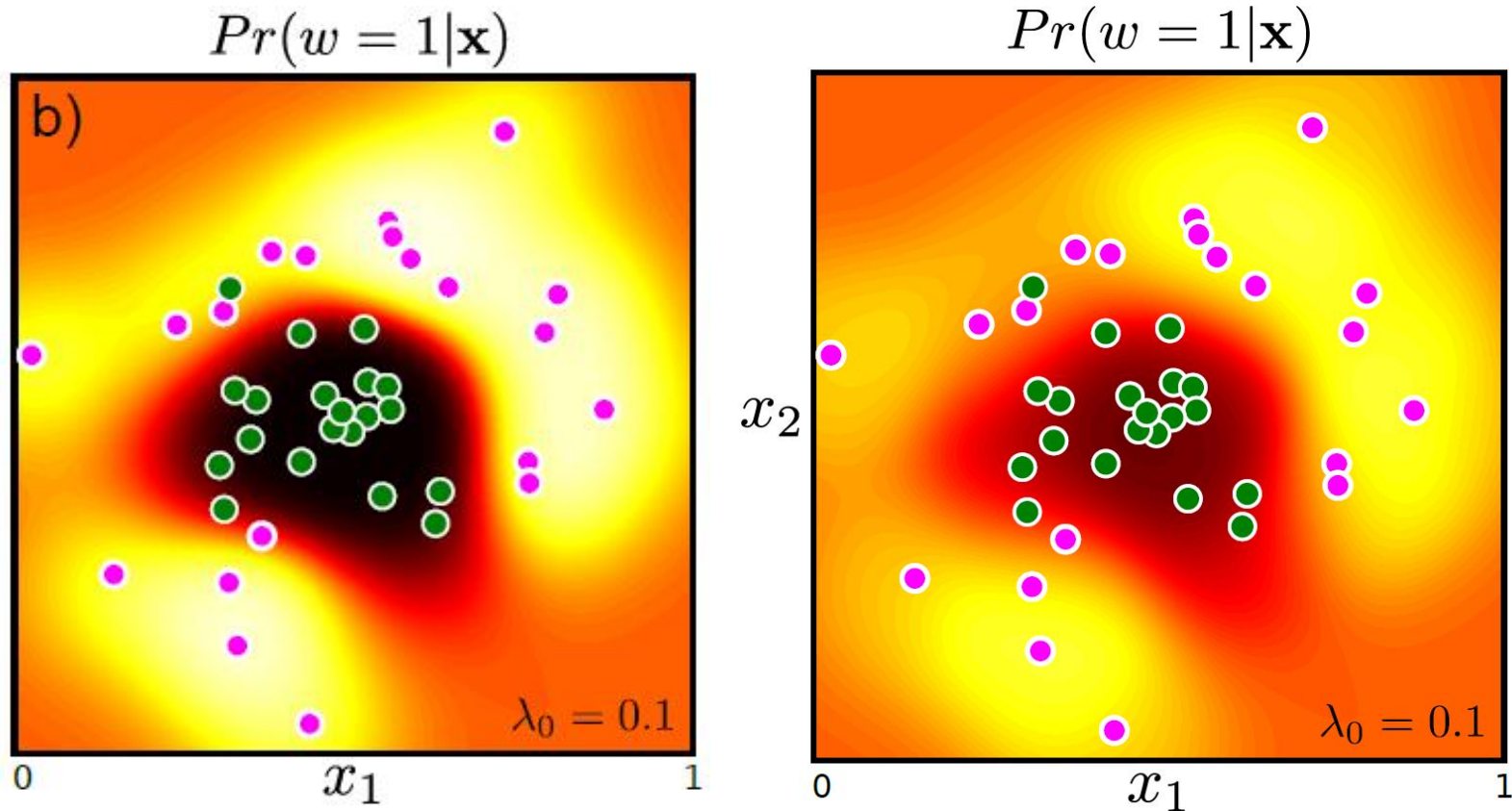
Depend only depend on inner products!

# Kernel Logistic Regression

$$k[\mathbf{x}_i, \mathbf{x}_j] = \exp \left[ -0.5 \left( \frac{(\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_j)}{\lambda_0^2} \right) \right]$$



# ML vs. Bayesian



Bayesian case is known as Gaussian process classification

# Relevance vector classification

Apply sparse prior to dual variables (dual logistic regression):

$$Pr(\boldsymbol{\psi}) = \prod_{i=1}^I \text{Stud}_{\psi_i} [0, 1, \nu]$$

As before, write as marginalization of dual variables:

$$\begin{aligned} Pr(\boldsymbol{\psi}) &= \prod_{i=1}^I \int \text{Norm}_{\psi_i} \left[ 0, \frac{1}{h_i} \right] \text{Gam}_{h_i} \left[ \frac{\nu}{2}, \frac{\nu}{2} \right] dh_i \\ &= \int \text{Norm}_{\boldsymbol{\psi}} [0, \mathbf{H}^{-1}] \prod_{i=1}^I \text{Gam}_{h_i} [\nu/2, \nu/2] d\mathbf{H}. \end{aligned}$$

# Relevance vector classification

Apply sparse prior to dual variables:

$$Pr(\boldsymbol{\psi}) = \int \text{Norm}_{\boldsymbol{\psi}}[0, \mathbf{H}^{-1}] \prod_{i=1}^I \text{Gam}_{h_i}[\nu/2, \nu/2] d\mathbf{H}.$$

Gives likelihood:

$$\begin{aligned} Pr(\mathbf{w}|\mathbf{X}) &= \int Pr(\mathbf{w}|\mathbf{X}, \boldsymbol{\psi}) Pr(\boldsymbol{\psi}) d\boldsymbol{\psi} \\ &= \iint \prod_{i=1}^I \text{Bern}_{w_i}[\text{sig}[\boldsymbol{\psi}^T \mathbf{K}[\mathbf{X}, \mathbf{x}_i]]] \text{Norm}_{\boldsymbol{\psi}}[0, \mathbf{H}^{-1}] \text{Gam}_{h_i}[\nu/2, \nu/2] d\mathbf{H} d\boldsymbol{\psi}. \end{aligned}$$

# Relevance vector classification

Laplace approximation:

$$Pr(\mathbf{w}|\mathbf{X}) \approx$$

$$\int \prod_{i=1}^I (2\pi)^{I/2} |\Sigma|^{0.5} \text{Bern}_{w_i}[\text{sig}[\boldsymbol{\mu}^T \mathbf{K}[\mathbf{X}, \mathbf{x}_i]]] \text{Norm}_{\boldsymbol{\mu}}[0, \mathbf{H}^{-1}] \text{Gam}_{h_i}\left[\frac{\nu}{2}, \frac{\nu}{2}\right] d\mathbf{H}$$

Second approximation:

$$Pr(\mathbf{w}|\mathbf{X}) \approx$$

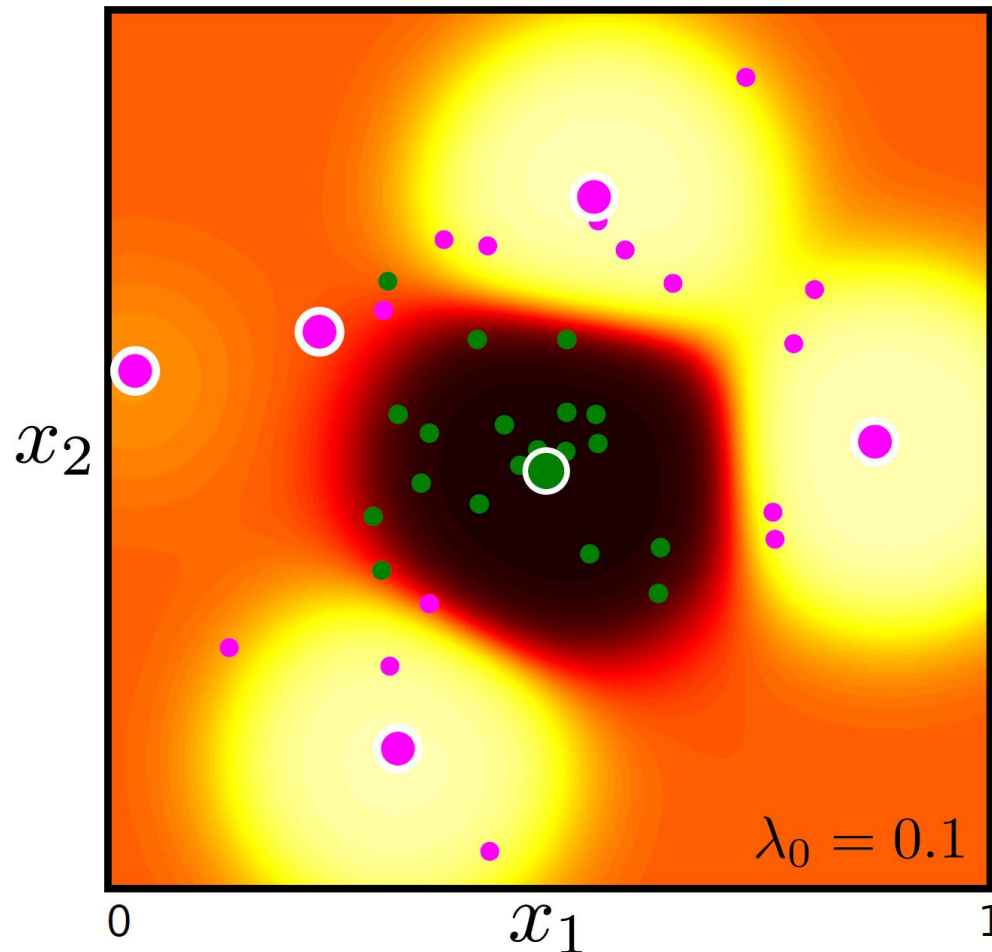
$$\max_{\mathbf{H}} \left[ \prod_{i=1}^I (2\pi)^{I/2} |\Sigma|^{0.5} \text{Bern}_{w_i}[\text{sig}[\boldsymbol{\mu}^T \mathbf{K}[\mathbf{X}, \mathbf{x}_i]]] \text{Norm}_{\boldsymbol{\mu}}[0, \mathbf{H}^{-1}] \text{Gam}_{h_i}\left[\frac{\nu}{2}, \frac{\nu}{2}\right] \right]$$

To solve, alternately update hidden variables in  $\mathbf{H}$  and mean and variance of Laplace approximation.



# Relevance vector classification

$$Pr(w = 1 | \mathbf{x})$$



The final solution only depends on a very small number of examples – efficient

The logo of the University of Bonn, featuring a blue square with a white curved line and a grey square.

UNIVERSITÄT **BONN**