



Machine Learning II: (Applications)

Prepared by:
Prof. Dr. Visvanathan Ramesh

References and Sources: Nils Bertschinger (ML II lecture slides)
Simon Prince (Learning and Vision)



- **Recap – Nil Bertschinger Lectures**
 - Approaches to Machine Learning
 - Introduction to Probability, Bayes rule, Probability Distributions
 - Bayesian Machine Learning
 - Parametric/Non-parametric Bayesian methods
 - Gaussian Processes
 - Link to Neural Networks
 - Bayesian Nonparametrics
 - Dirichlet Processes, Chinese Restaurant Process, Indian Buffet Process
 - Inference by Sampling , MCMC – Metropolis-Hastings, Gibbs Sampler, HMC sampler
- What is not covered yet? (And what is planned for the rest of the weeks)
 - Recap of past classes , discussion
 - Example application case studies in computer vision
 - Variational Methods



Brief Recap



- ▶ Data-driven
 - ▶ Very large data sets ... “Big Data”
 - ▶ Non-parametric models, e.g. k-NN
- ▶ Model-driven
 - ▶ Can be used for small data sets
 - ▶ Parametric models

Note: As models become more complex any data set is “small”
⇒ Recent rise of model based machine learning

General setup of model based ML:

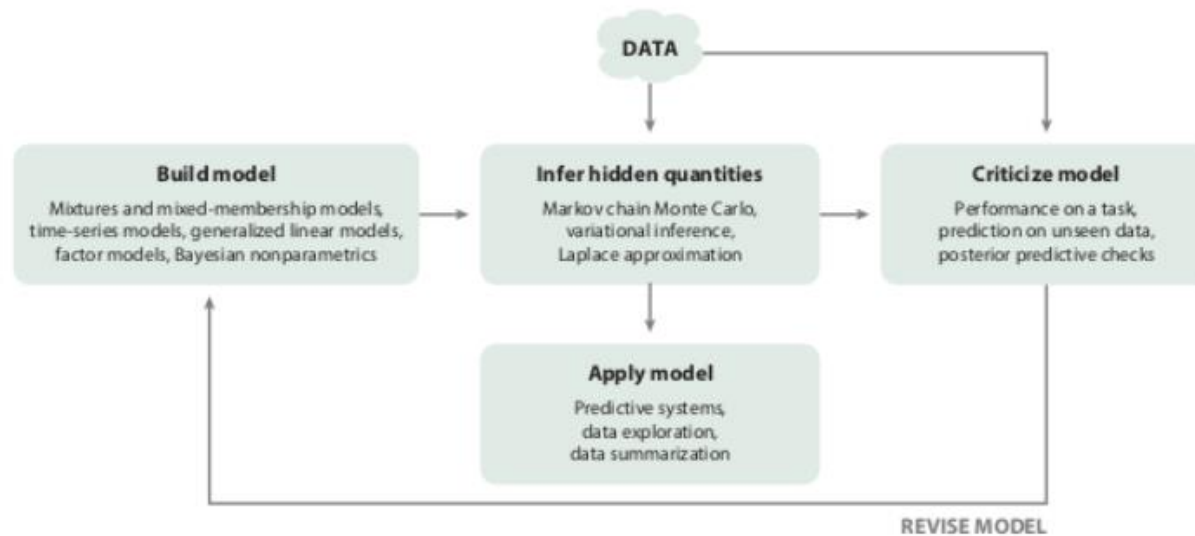


Fig. from: David M. Blei, *Build, Compute, Critique, Repeat: Data Analysis with Latent Variable Models*, Annu.

Rev. Stat. Appl. 2014. 1:20332



- ▶ **Supervised:** Patterns whose class/output is known a-priori are used for training (*labelled training data*)
 - ▶ *Regression:* Real-valued output
Typical examples: Interpolation, (Time-series) Prediction
 - ▶ *Classification:* Categorical output
Typical examples: Face recognition, Identity authentication, Speech recognition
- ▶ **Unsupervised:** Number of classes is (in general) unknown and no labelled data are available
Typical examples: Cluster analysis, Recommendation systems

Bayesian statistics:

- ▶ Principled and logically consistent way to reason under uncertainty

Prior $\xrightarrow{\text{Data}}$ Posterior (belief update)

- ▶ Especially useful when taking decisions or making predictions

Bayesian machine learning:

- ▶ Statistical modeling:

$$p(\underbrace{\mathbf{x}}_{\text{Data}}, \underbrace{\mathbf{z}}_{\substack{\text{Latent variables} \\ \text{Parameters}}}) = \underbrace{p(\mathbf{z})}_{\text{prior}} \underbrace{p(\mathbf{x}|\mathbf{z})}_{\text{likelihood}}$$

- ▶ Conceptually simple, but computationally challenging



Bayesian machine learning:

- ▶ Bayesian modeling requires prior assumptions:
 - ▶ Parametric models, e.g. linear regression
 - ▶ Bayesian non-parametrics:
 - ▶ Flexible models with infinite-dimensional parameter spaces
 - ▶ Effective number of parameters grow with amount of data

But, explicit about prior assumptions

- ▶ *No free lunch theorem*: Assumption-free learning is impossible!
- ▶ Takes uncertainty into account
Bayesian Occam's razor: Automatic penalty for model complexity
- ▶ Computational challenge: Posterior $p(\mathbf{z}|\mathbf{x})$ often intractable
 - ▶ Sampling algorithms
 - ▶ Variational approximations



Machine Learning II course ... Focus on Bayesian methods

- ▶ Motivation: Bayesian vs frequentist statistics
- ▶ Decision theory: Handling uncertainty, loss functions
- ▶ Probability theory: Conjugate priors
- ▶ Modeling: Latent variables, hierarchical models, Bayesian non-parametrics
- ▶ Model selection: Marginal likelihood, sparsity priors
- ▶ Algorithms: Variational Bayes (ELBO), sampling methods

Potential applications

- ▶ Social data: Voting results, network models
- ▶ Economic data: GDP forecasting, volatility modeling
- ▶ Computer vision: Detection, tracking, recognition, segmentation
- ▶ ...



Computer vision: models, learning and inference

Source: Chapter 6 , 7

**Computer Vision: Models, Learning and Inference
(Simon Prince)**

Structure



Computer vision models

- Two types of model

Worked example 1: Regression

Worked example 2: Classification

Which type should we choose?

Applications

Observe **measured data**, x

Draw inferences from it about **state of world**, w

Examples:

- Observe adjacent frames in video sequence
- Infer camera motion

- Observe image of face
- Infer identity

- Observe images from two displaced cameras
- Infer 3d structure of scene

Regression vs. Classification



Observe measured data, x

Draw inferences from it about world, w

When the world state w is **continuous** we'll call this
regression

When the world state w is **discrete** we call this
classification



Unfortunately visual measurements may be compatible with more than one world state w

- Measurement process is noisy
- Inherent ambiguity in visual data

Conclusion: the best we can do is compute a probability distribution $\Pr(w|x)$ over possible states of world

Refined goal of computer vision



Take observations x

Return probability distribution $\Pr(w|x)$ over possible worlds compatible with data

(not always tractable – might have to settle for an approximation to this distribution, samples from it, or the best (MAP) solution for w)

Components of solution



We need

A **model** that mathematically relates the visual data x to the world state w . Model specifies family of relationships, particular relationship depends on parameters θ

A **learning algorithm**: fits parameters θ from paired training examples x_i, w_i

An **inference algorithm**: uses model to return $\Pr(w|x)$ given new observed data x .



The **model** mathematically relates the visual data x to the world state w . Two main categories of model

1. Model contingency of the world on the data $\Pr(w|x)$
2. Model contingency of data on world $\Pr(x|w)$



1. Model contingency of the world on the data
 $\Pr(w|x)$

(DISCRIMINATIVE MODEL)

2. Model contingency of data on world $\Pr(x|w)$

(GENERATIVE MODELS)

**Generative as probability model over data and
so when we draw samples from model, we
GENERATE new data**



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

1. Choose an appropriate form for $\Pr(w)$
2. Make parameters a function of x
3. Function takes parameters θ that define its shape

Learning algorithm: learn parameters θ from training data x, w

Inference algorithm: just evaluate $\Pr(w|x)$

Type 2: $\Pr(\mathbf{x}|\mathbf{w})$ - Generative



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

1. Choose an appropriate form for $\Pr(\mathbf{x})$
2. Make parameters a function of \mathbf{w}
3. Function takes parameters θ that define its shape

Learning algorithm: learn parameters θ from training data \mathbf{x}, \mathbf{w}

Inference algorithm: Define prior $\Pr(\mathbf{w})$ and then compute $\Pr(\mathbf{w}|\mathbf{x})$ using Bayes' rule

$$\Pr(\mathbf{w}|\mathbf{x}) = \frac{\Pr(\mathbf{x}|\mathbf{w})\Pr(\mathbf{w})}{\int \Pr(\mathbf{x}|\mathbf{w})\Pr(\mathbf{w})d\mathbf{w}}$$



Two different types of model depend on the quantity of interest:

- 1. $\Pr(w|x)$ Discriminative**
- 2. $\Pr(w|x)$ Generative**

**Inference in discriminative models easy as we directly model posterior $\Pr(w|x)$.
Generative models require more complex inference process using Bayes' rule**

Structure



Computer vision models

- Two types of model

Worked example 1: Regression

Worked example 2: Classification

Which type should we choose?

Applications

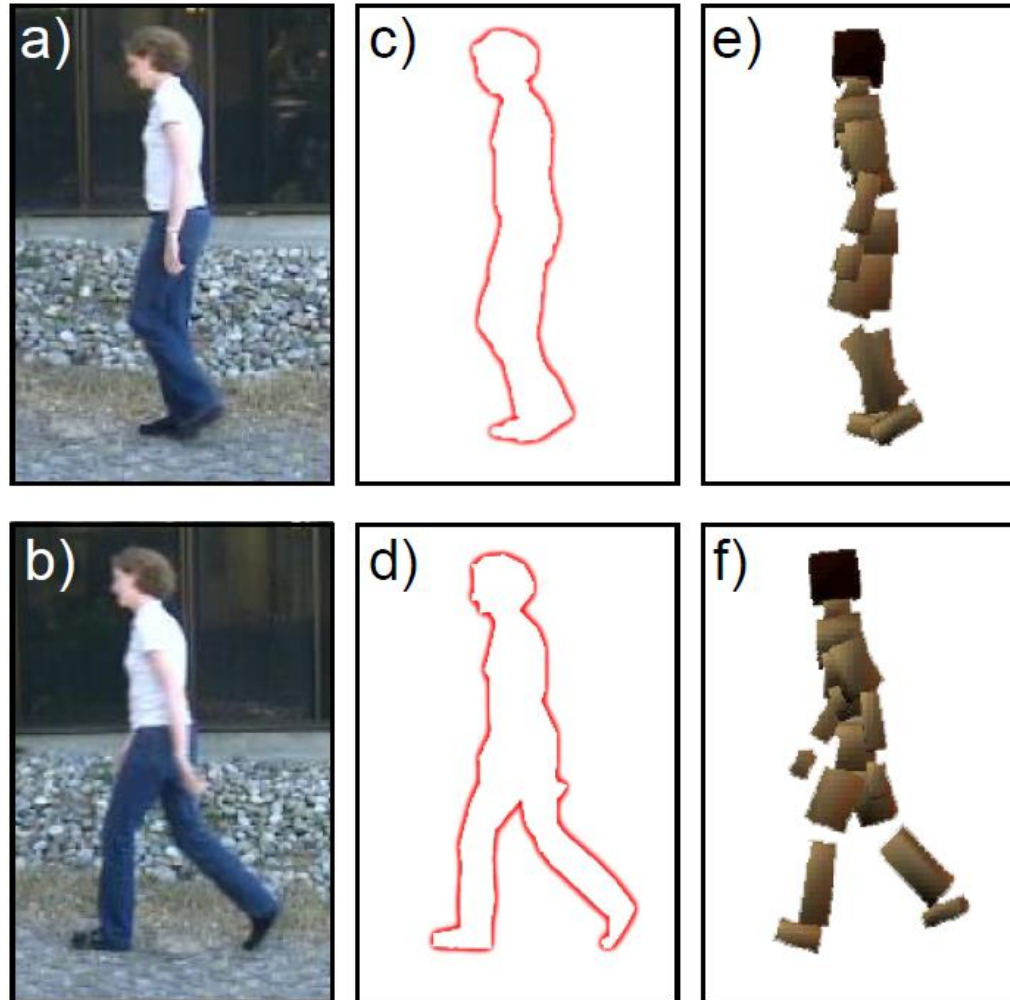
Worked example 1: Regression

Consider simple case where

- we make a univariate continuous measurement \mathbf{x}
- use this to predict a univariate continuous state \mathbf{w}

(regression as world state is continuous)

Regression application 1: Pose from Silhouette



Regression application 2: Head pose estimation



-76°



-11°



2°



8°



43°



79°

Worked example 1: Regression

Consider simple case where

- we make a univariate continuous measurement \mathbf{x}
- use this to predict a univariate continuous state \mathbf{w}

(regression as world state is continuous)



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

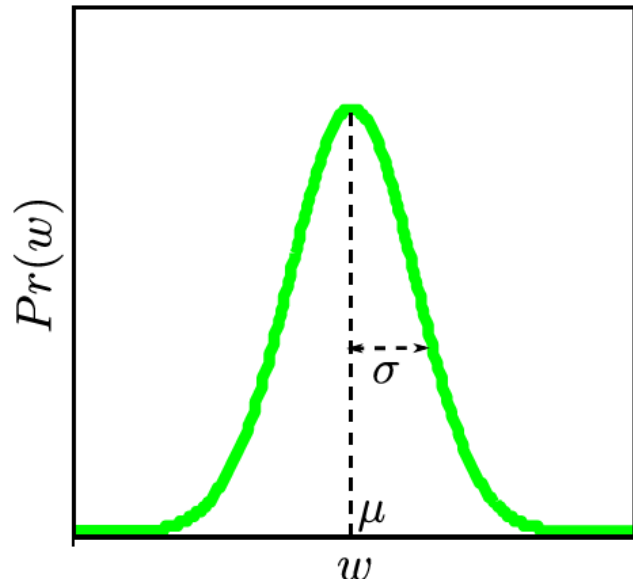
1. Choose an appropriate form for $\Pr(w)$
2. Make parameters a function of \mathbf{x}
3. Function takes parameters θ that define its shape

Learning algorithm: learn parameters θ from training data \mathbf{x}, w

Inference algorithm: just evaluate $\Pr(w|x)$

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

1. Choose an appropriate form for $\Pr(w)$
2. Make parameters a function of \mathbf{x}
3. Function takes parameters θ that define its shape

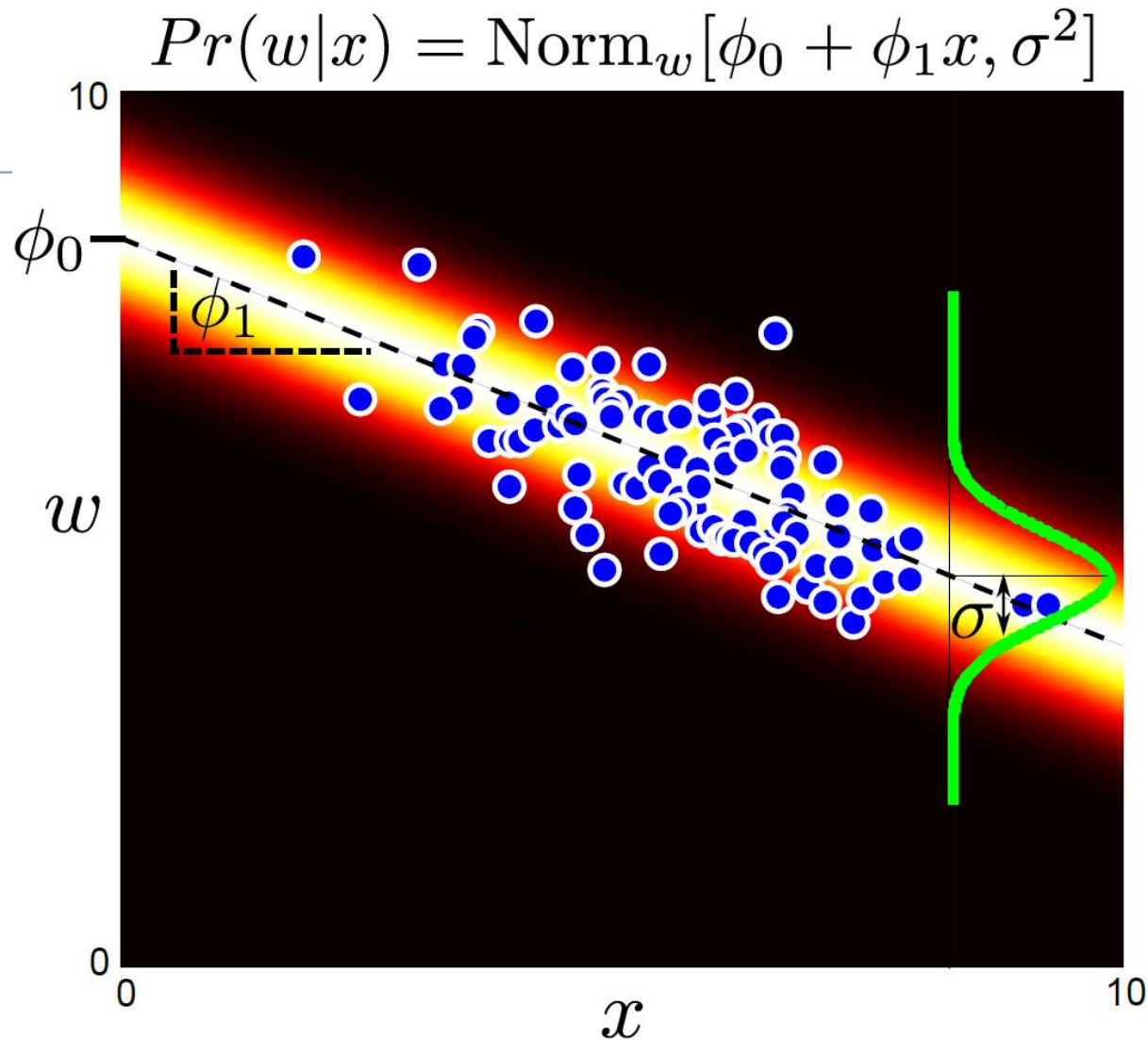


1. Choose normal distribution over w
2. Make mean μ linear function of x
(variance constant)

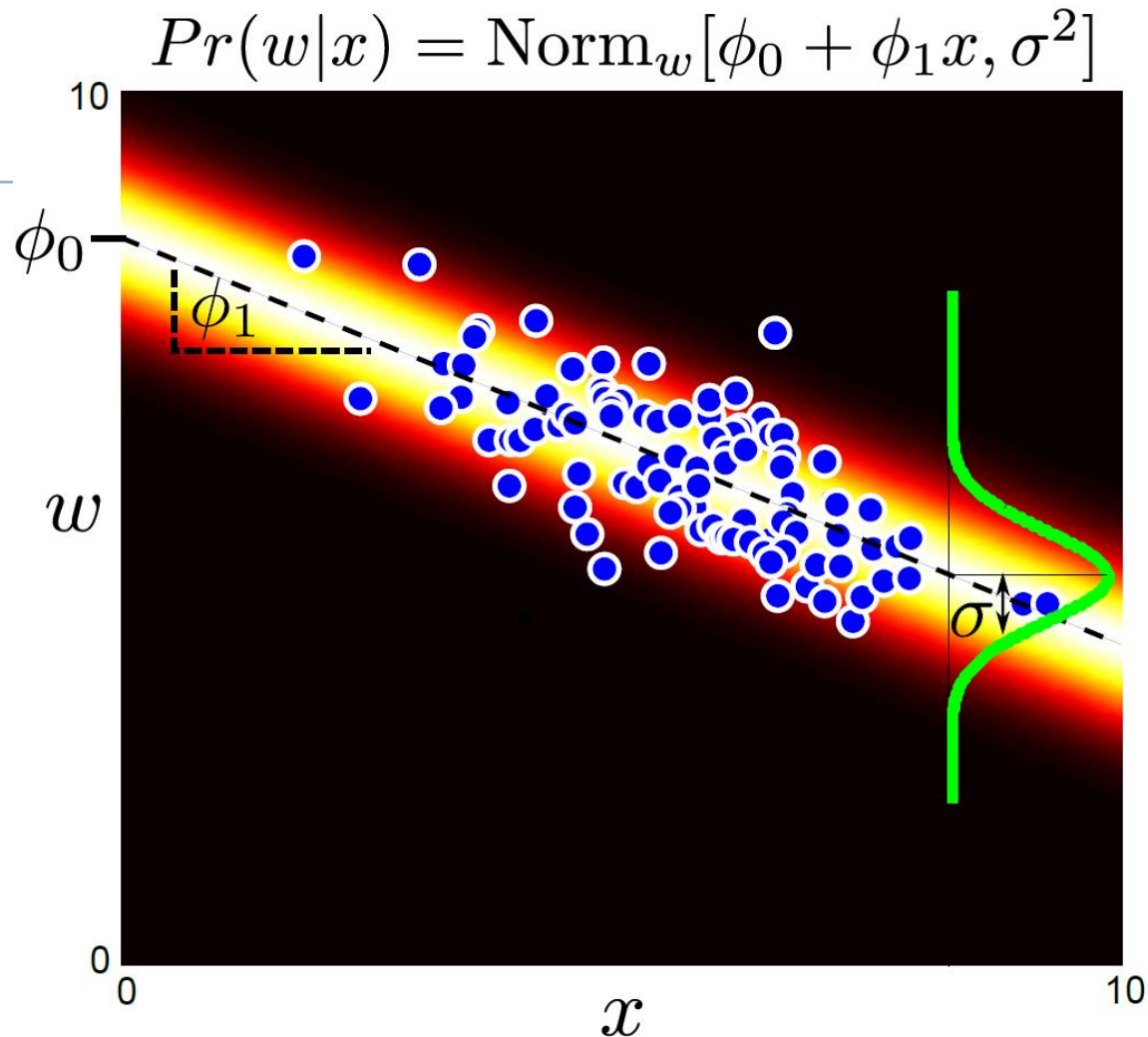
$$\Pr(w|x, \theta) = \text{Norm}_w [\phi_0 + \phi_1 x, \sigma^2]$$

3. Parameters are ϕ_0, ϕ_1, σ^2 .

This model is called *linear regression*.

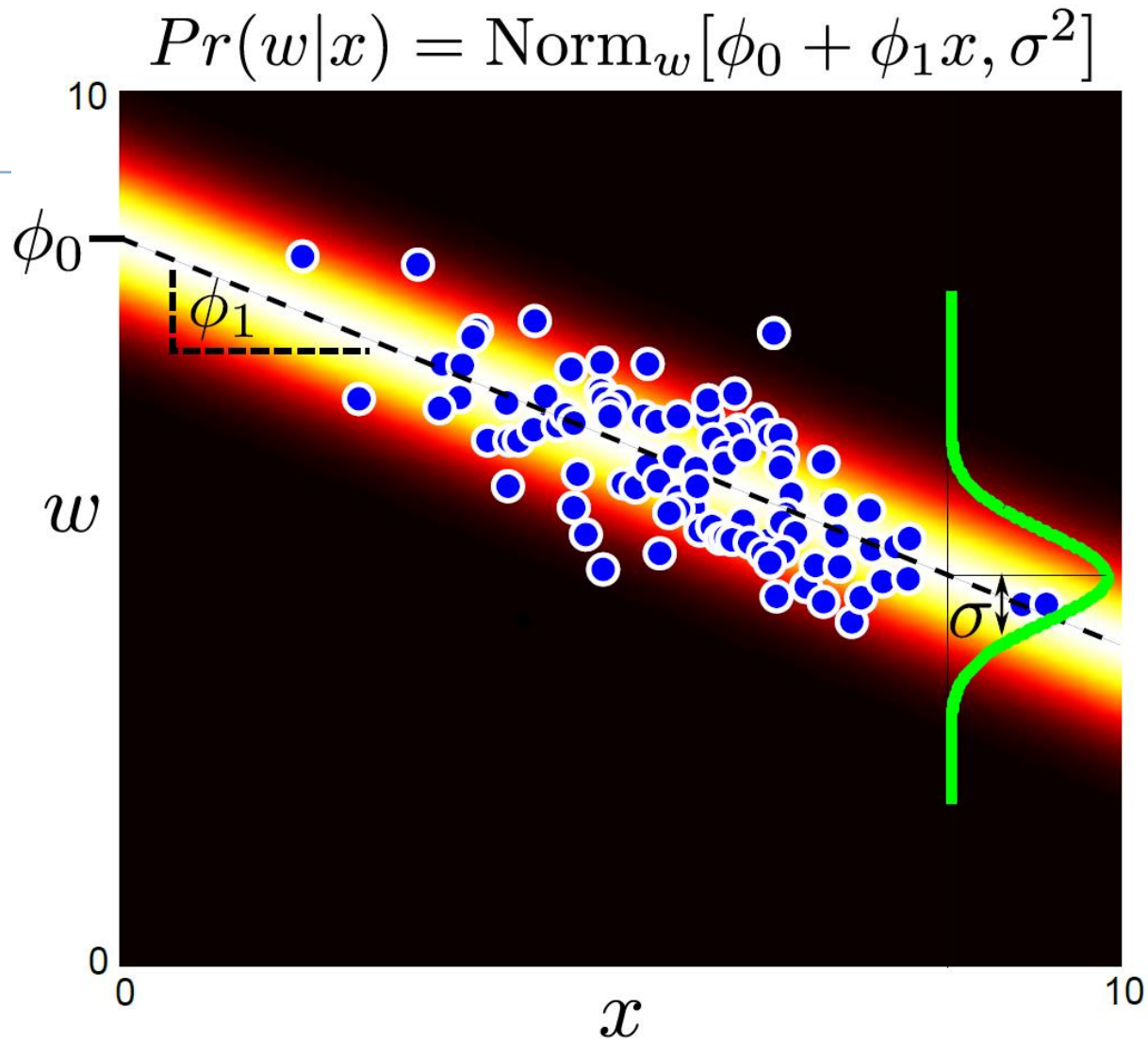


Parameters $\theta = \{\phi_0, \phi_1, \sigma^2\}$ are y-offset, slope and variance



Learning algorithm: learn θ from training data \mathbf{x}, \mathbf{y} . E.g.

$$\begin{aligned} \hat{\theta} &= \arg \max_{\theta} Pr(\theta | w_{1...I}, x_{1...I}) \\ &= \arg \max_{\theta} Pr(w_{1...I} | x_{1...I}, \theta) Pr(\theta) &= \arg \max_{\theta} \prod_{i=1}^I Pr(w_i | x_i, \theta) Pr(\theta), \end{aligned}$$



Inference algorithm: just evaluate $Pr(w|x)$ for new data

Type 2: $\Pr(\mathbf{x}|\mathbf{w})$ - Generative



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

1. Choose an appropriate form for $\Pr(\mathbf{x})$
2. Make parameters a function of \mathbf{w}
3. Function takes parameters θ that define its shape

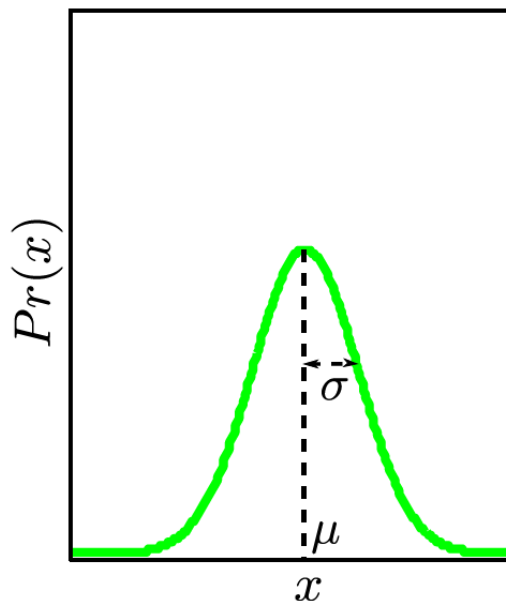
Learning algorithm: learn parameters θ from training data \mathbf{x}, \mathbf{w}

Inference algorithm: Define prior $\Pr(\mathbf{w})$ and then compute $\Pr(\mathbf{w}|\mathbf{x})$ using Bayes' rule

$$\Pr(\mathbf{w}|\mathbf{x}) = \frac{\Pr(\mathbf{x}|\mathbf{w})\Pr(\mathbf{w})}{\int \Pr(\mathbf{x}|\mathbf{w})\Pr(\mathbf{w})d\mathbf{w}}$$

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

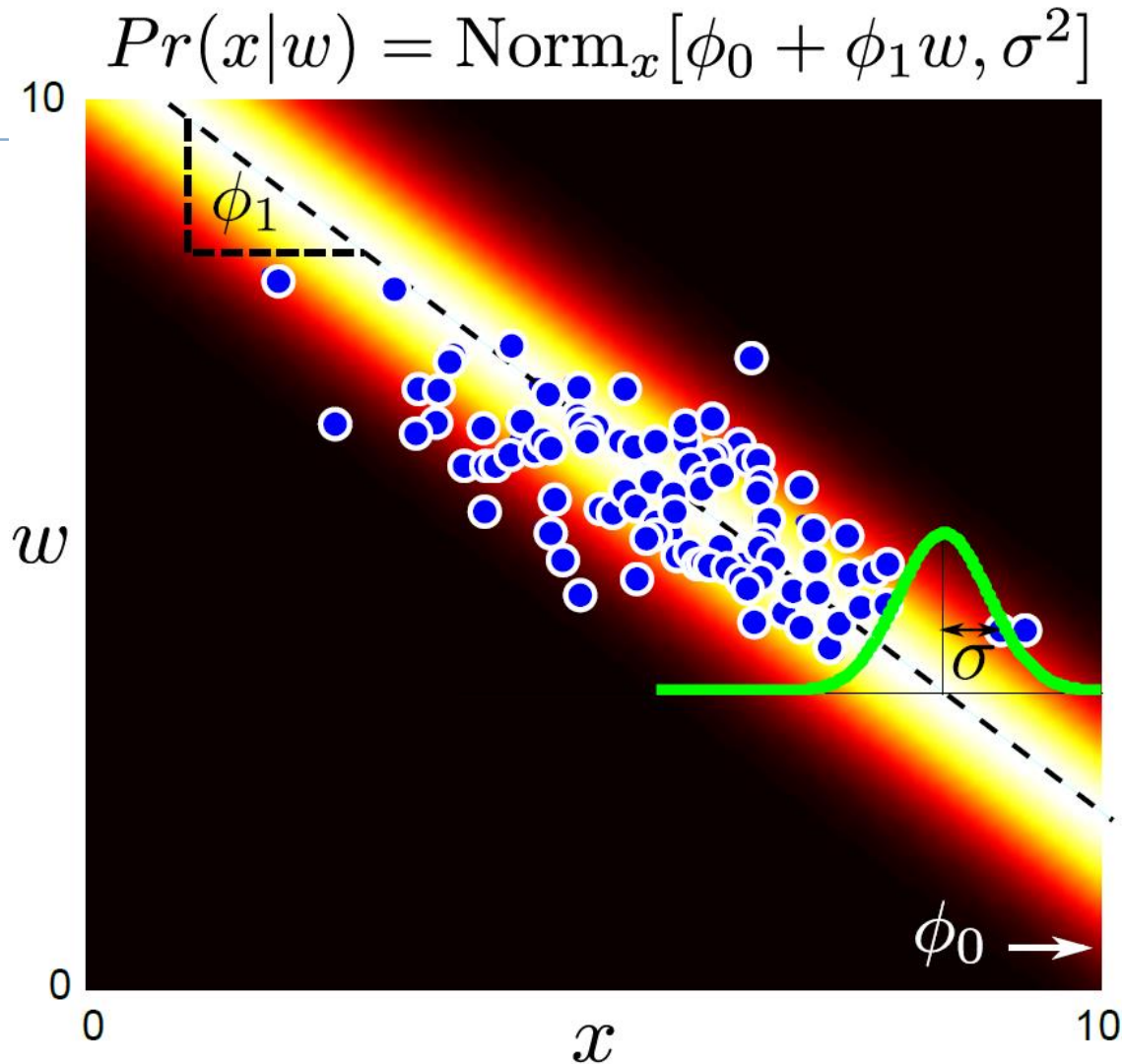
1. Choose an appropriate form for $\Pr(\mathbf{x})$
2. Make parameters a function of \mathbf{w}
3. Function takes parameters θ that define its shape



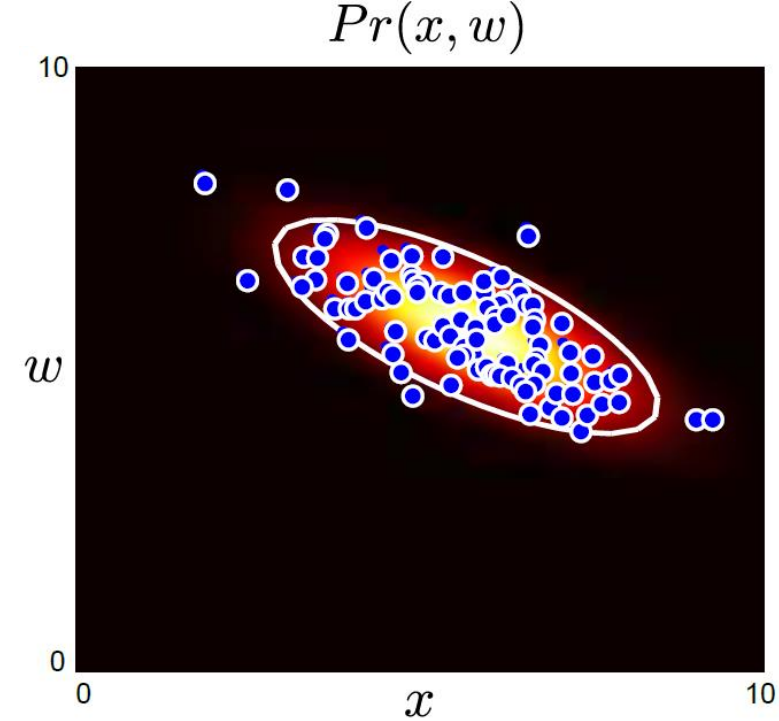
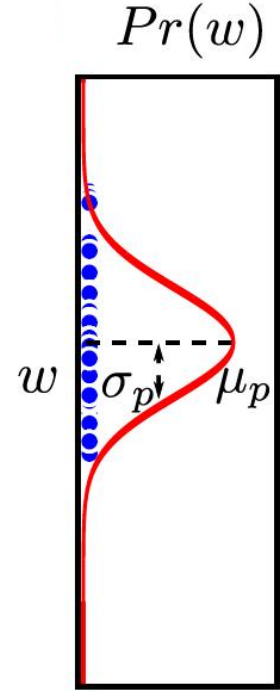
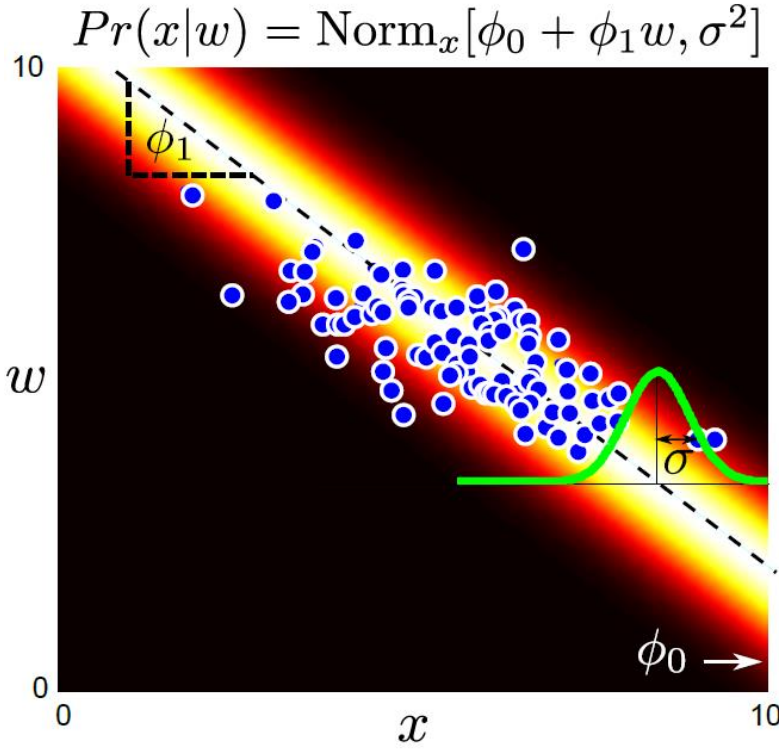
1. Choose normal distribution over x
2. Make mean μ linear function of w
(variance constant)

$$\Pr(x|w, \theta) = \text{Norm}_x [\phi_0 + \phi_1 w, \sigma^2]$$

3. Parameter are ϕ_0, ϕ_1, σ^2 .

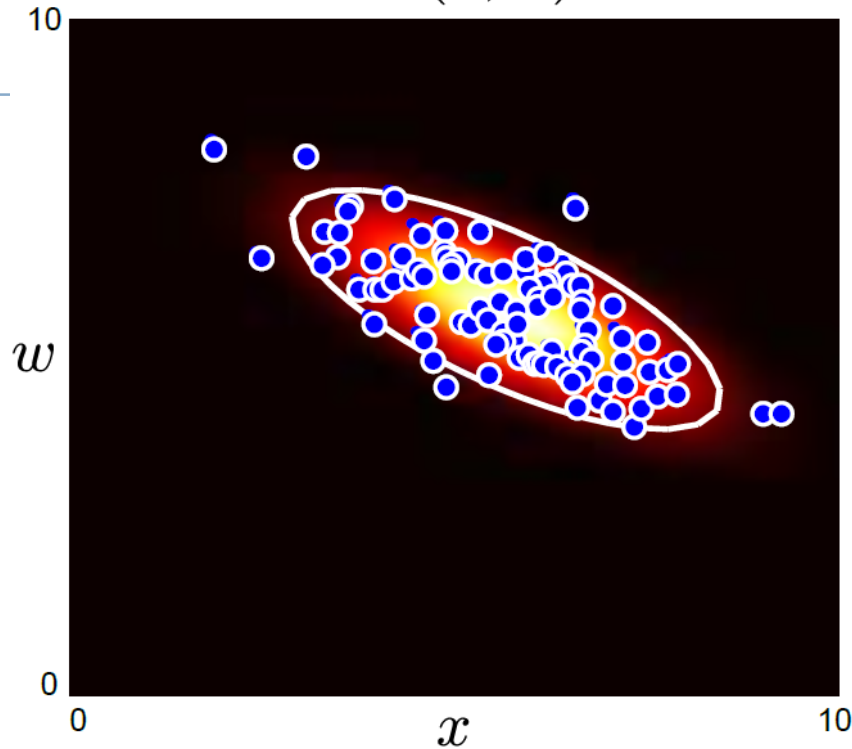
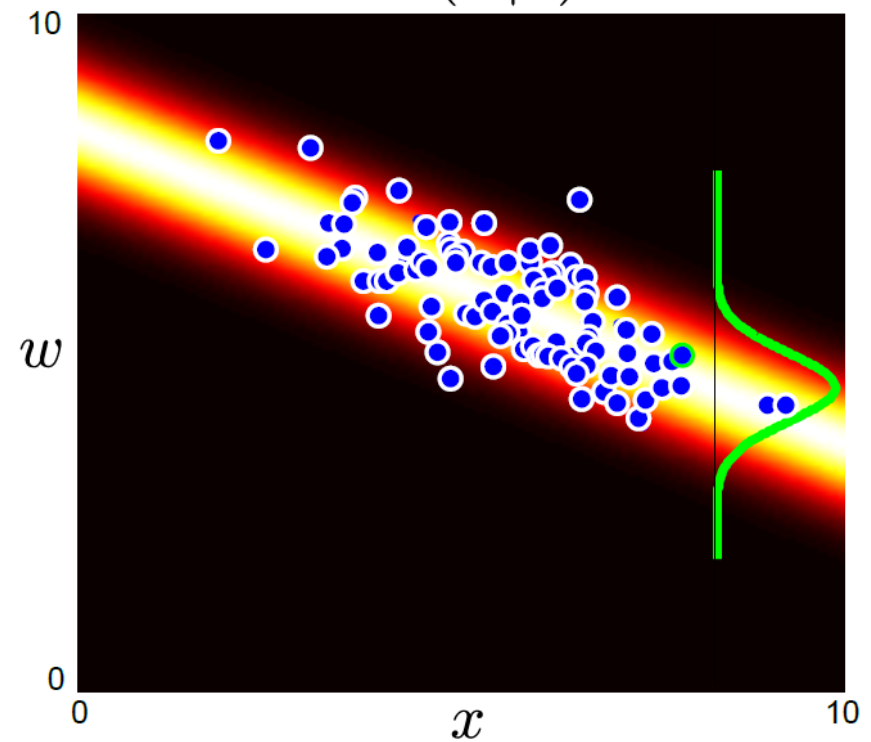


Learning algorithm: learn θ from training data \mathbf{x}, \mathbf{w} . e.g.
MAP



$$Pr(x|w) \quad x \quad Pr(w) \quad = \quad Pr(x, w)$$

Can get back to joint probability $Pr(x, y)$

$Pr(x, w)$  $Pr(w|x)$ 

Inference algorithm: compute $Pr(\mathbf{w}|\mathbf{x})$ using Bayes rule

$$Pr(\mathbf{w}|\mathbf{x}) = \frac{Pr(\mathbf{x}|\mathbf{w})Pr(\mathbf{w})}{\int Pr(\mathbf{x}|\mathbf{w})Pr(\mathbf{w})d\mathbf{w}}$$

Structure



Computer vision models

- Three types of model

Worked example 1: Regression

Worked example 2: Classification

Which type should we choose?

Applications

Worked example 2: Classification

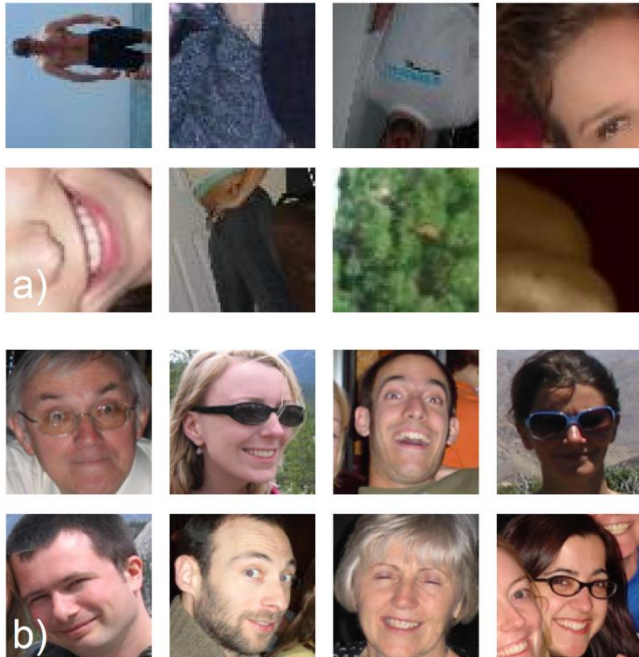


Consider simple case where

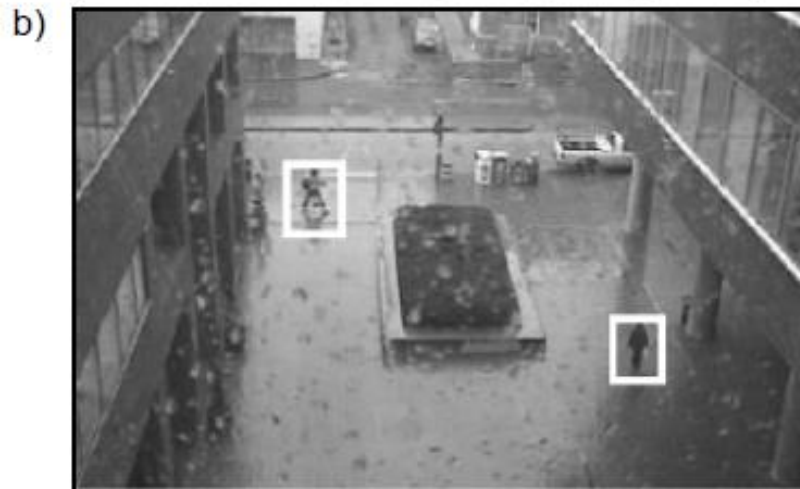
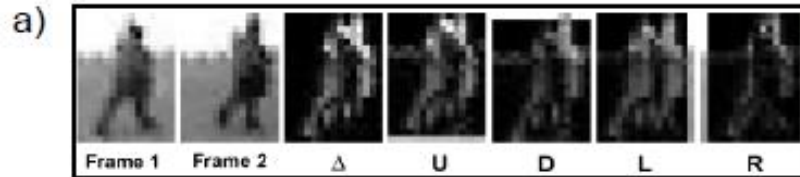
- we make a univariate continuous measurement x
- use this to predict a discrete binary $w \in \{0, 1\}$

(classification as world state is discrete)

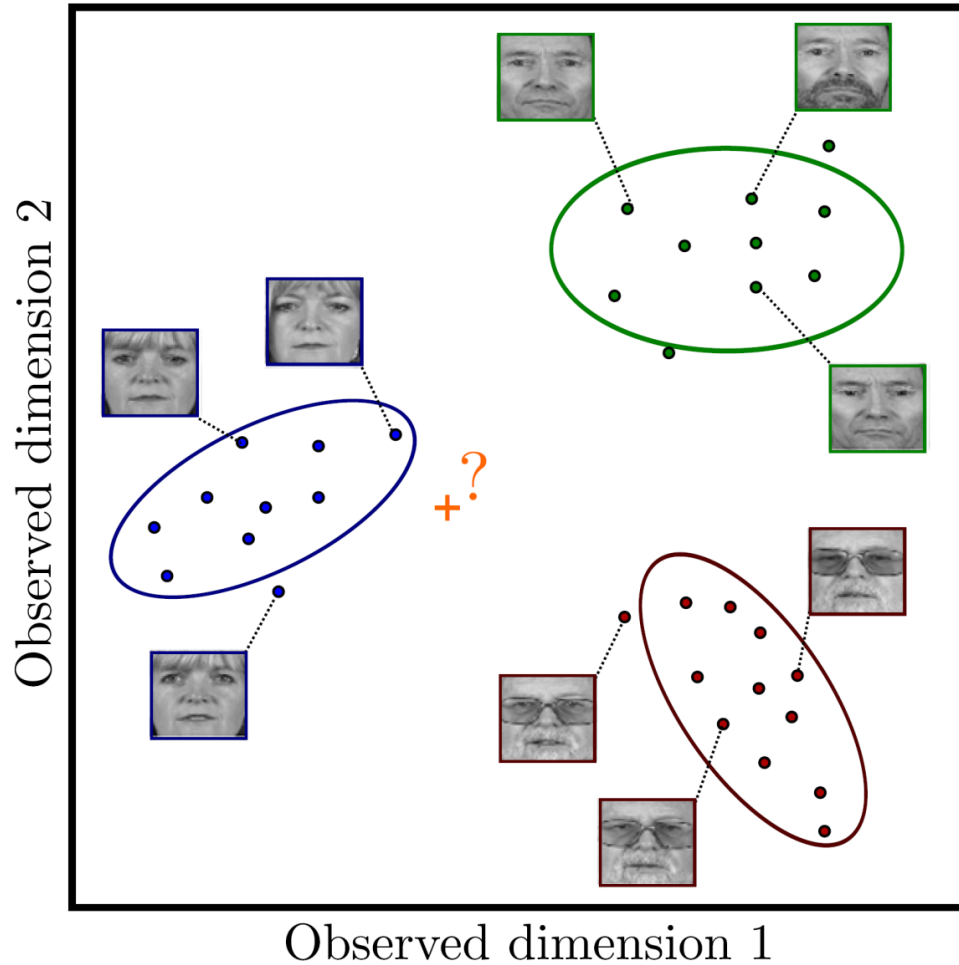
Classification Example 1: Face Detection



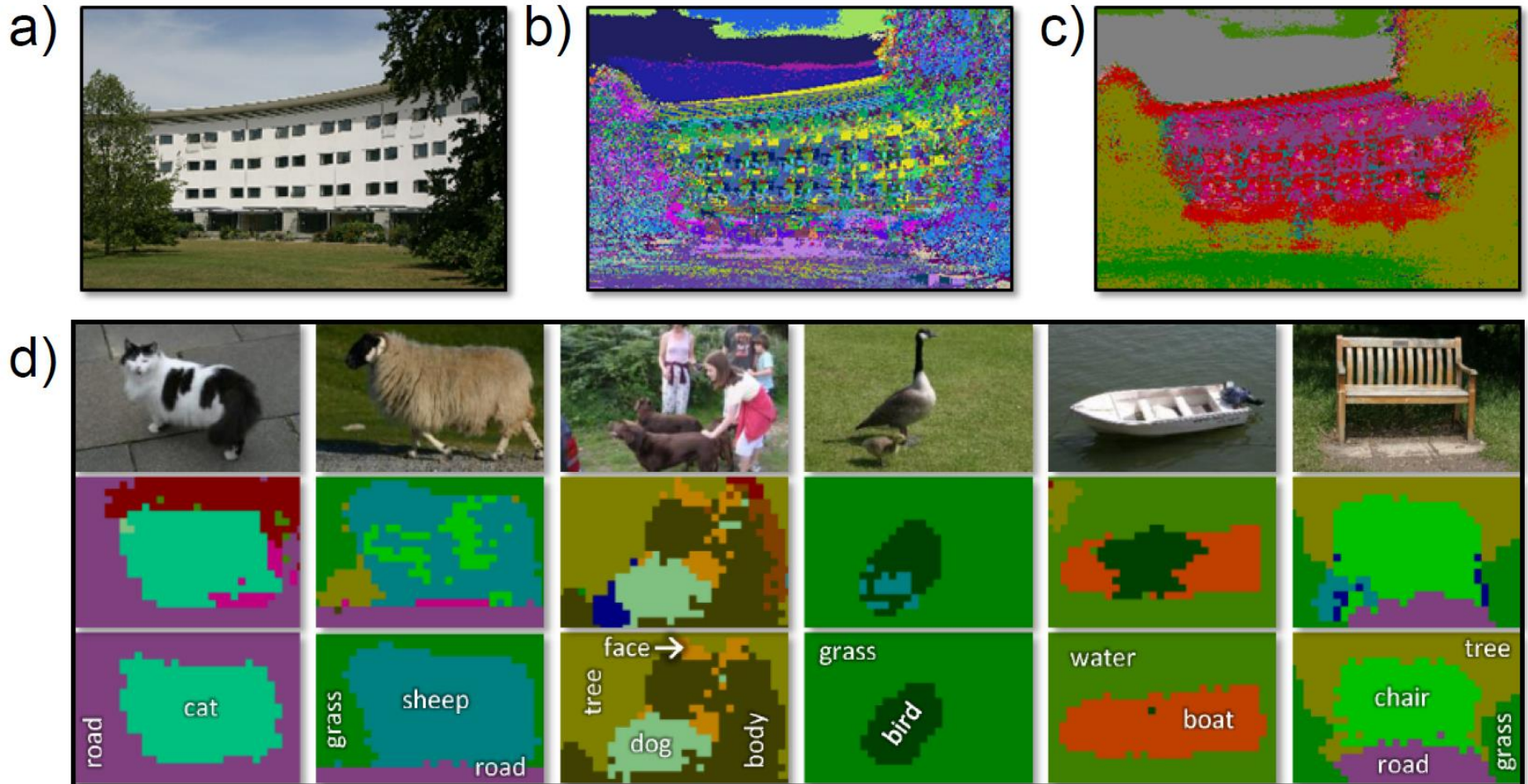
Classification Example 2: Pedestrian Detection



Classification Example 3: Face Recognition



Classification Example 4: Semantic Segmentation



Worked example 2: Classification



Consider simple case where

- we make a univariate continuous measurement x
- use this to predict a discrete binary world
 $w \in \{0, 1\}$

(classification as world state is discrete)



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

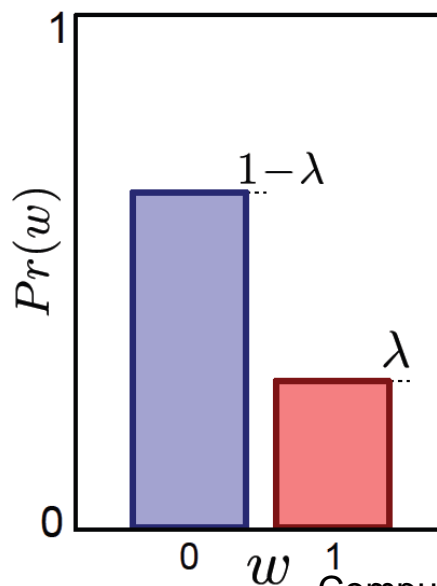
- Choose an appropriate form for $\Pr(w)$
- Make parameters a function of \mathbf{x}
- Function takes parameters θ that define its shape

Learning algorithm: learn parameters θ from training data \mathbf{x}, w

Inference algorithm: just evaluate $\Pr(w|x)$

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

1. Choose an appropriate form for $\Pr(\mathbf{w})$
2. Make parameters a function of \mathbf{x}
3. Function takes parameters θ that define its shape

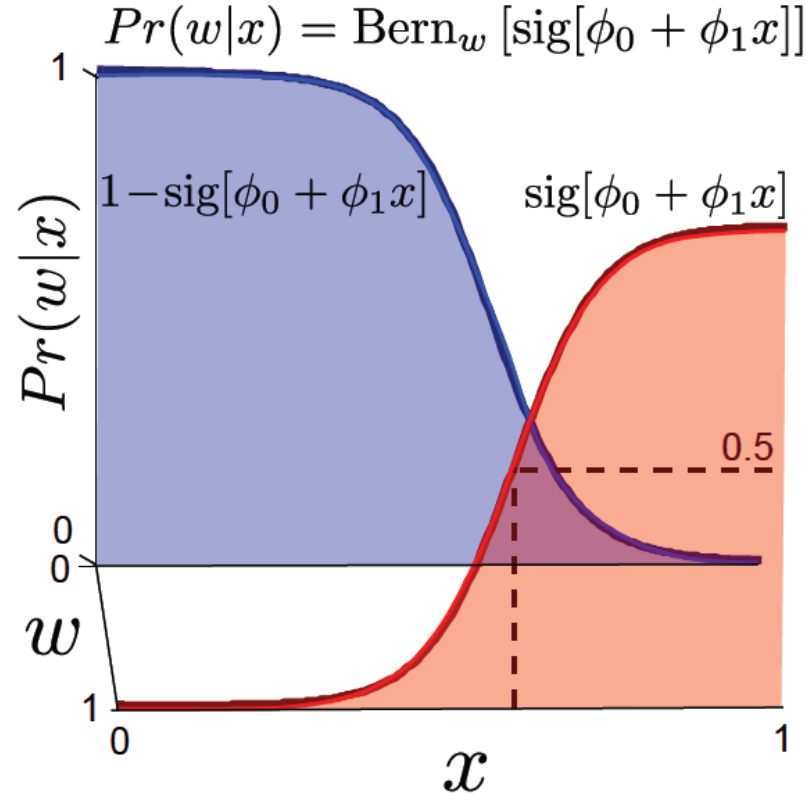
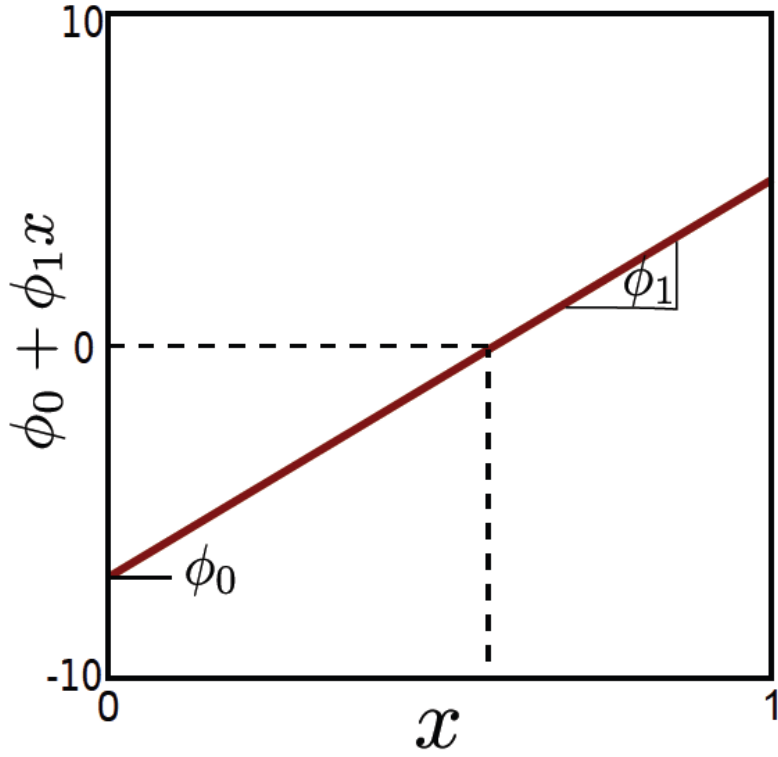


1. Choose Bernoulli dist. for $\Pr(\mathbf{w})$
2. Make parameters a function of \mathbf{x}

$$\Pr(w|x) = \text{Bern}_w [\text{sig}[\phi_0 + \phi_1 x]]$$

$$= \text{Bern}_w \left[\frac{1}{1 + \exp[-\phi_0 - \phi_1 x]} \right]$$

3. Function takes parameters ϕ_0 and ϕ_1
This model is called *logistic regression*.



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

Learning by standard methods (ML, MAP, Bayesian)

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

Type 2: $\Pr(\mathbf{x}|\mathbf{w})$ - Generative



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

1. Choose an appropriate form for $\Pr(\mathbf{x})$
2. Make parameters a function of \mathbf{w}
3. Function takes parameters θ that define its shape

Learning algorithm: learn parameters θ from training data \mathbf{x}, \mathbf{w}

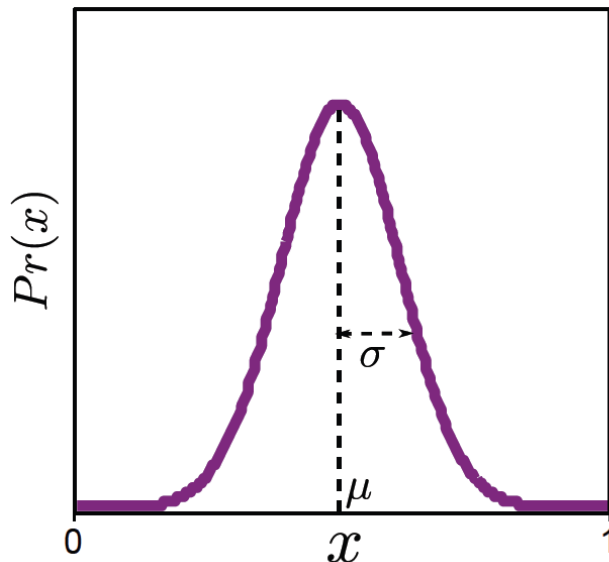
Inference algorithm: Define prior $\Pr(\mathbf{w})$ and then compute $\Pr(\mathbf{w}|\mathbf{x})$ using Bayes' rule

$$\Pr(\mathbf{w}|\mathbf{x}) = \frac{\Pr(\mathbf{x}|\mathbf{w})\Pr(\mathbf{w})}{\int \Pr(\mathbf{x}|\mathbf{w})\Pr(\mathbf{w})d\mathbf{w}}$$

Type 2: $\Pr(\mathbf{x}|\mathbf{w})$ - Generative

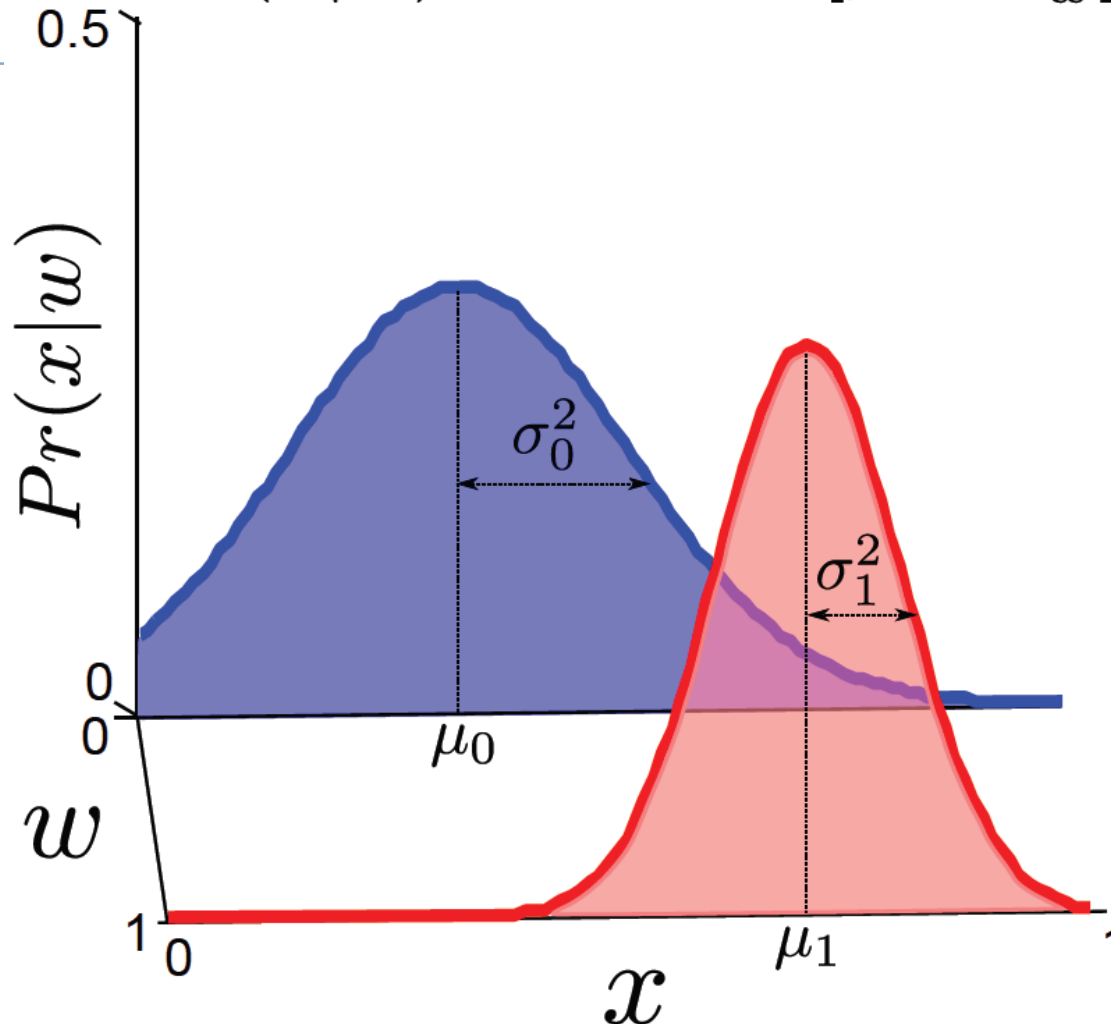
How to model $\Pr(\mathbf{x}|\mathbf{w})$?

1. Choose an appropriate form for $\Pr(\mathbf{x})$
2. Make parameters a function of \mathbf{w}
3. Function takes parameters θ that define its shape



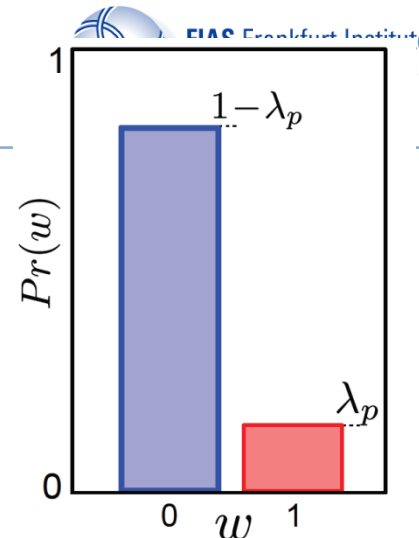
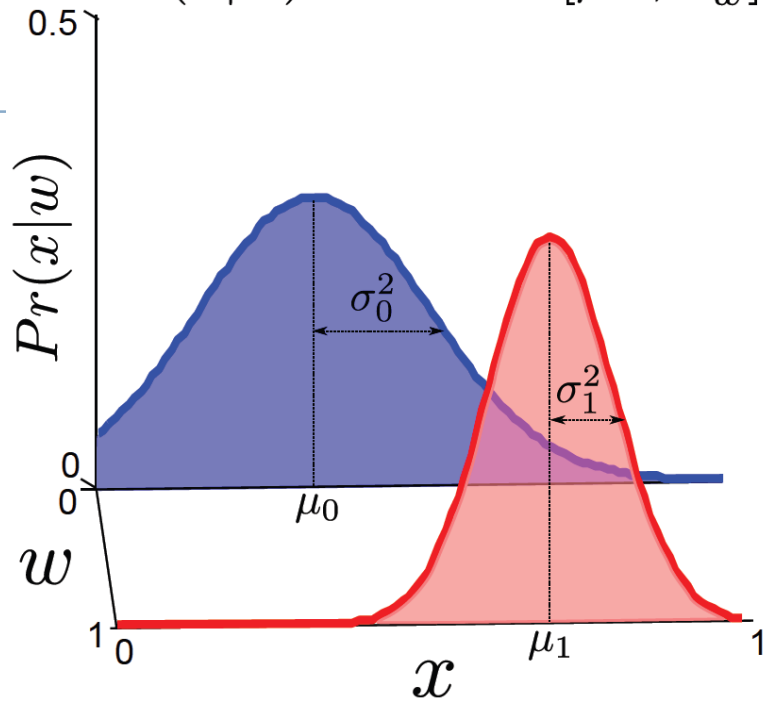
1. Choose a Gaussian distribution for $\Pr(\mathbf{x})$
2. Make parameters a function of discrete binary \mathbf{w}
$$\Pr(x|\mathbf{w}) = \text{Norm}_x[\mu_{\mathbf{w}}, \sigma_{\mathbf{w}}^2]$$
3. Function takes parameters $\mu_0, \mu_1, \sigma_0^2, \sigma_1^2$ that define its shape

$$Pr(x|w) = \text{Norm}_x[\mu_w, \sigma_w^2]$$



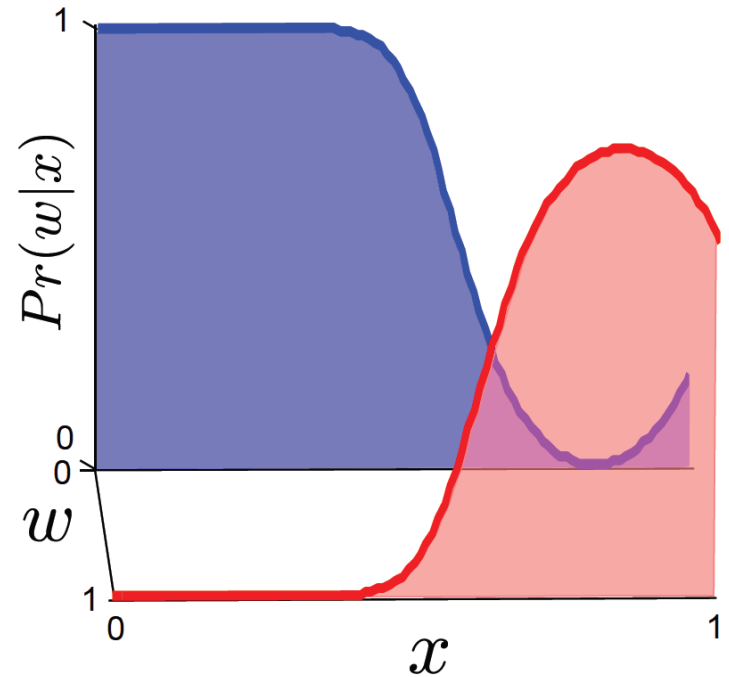
Learn parameters $\mu_0, \mu_1, \sigma_0^2, \sigma_1^2$ that define its shape

$$Pr(x|w) = \text{Norm}_x[\mu_w, \sigma_w^2]$$



Inference algorithm: Define prior $Pr(\mathbf{w})$ and then compute $Pr(\mathbf{w}|\mathbf{x})$ using Bayes' rule

$$Pr(\mathbf{w}|\mathbf{x}) = \frac{Pr(\mathbf{x}|\mathbf{w})Pr(\mathbf{w})}{\int Pr(\mathbf{x}|\mathbf{w})Pr(\mathbf{w})d\mathbf{w}}$$



Structure



Computer vision models

- Three types of model

Worked example 1: Regression

Worked example 2: Classification

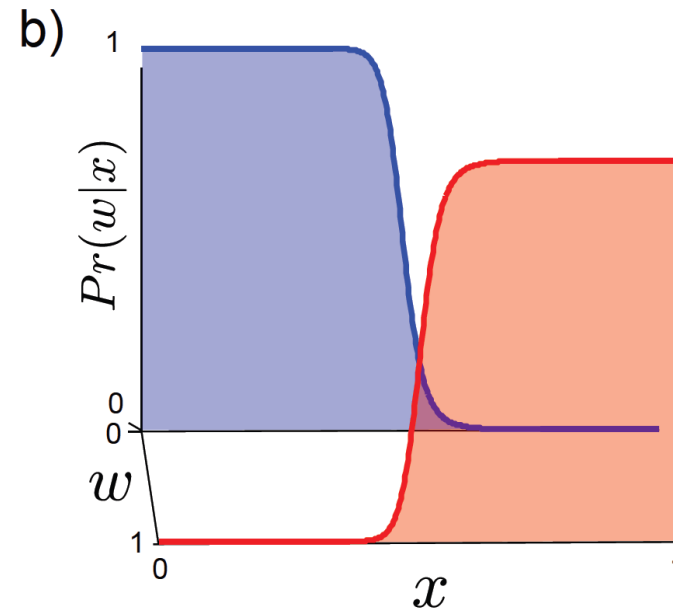
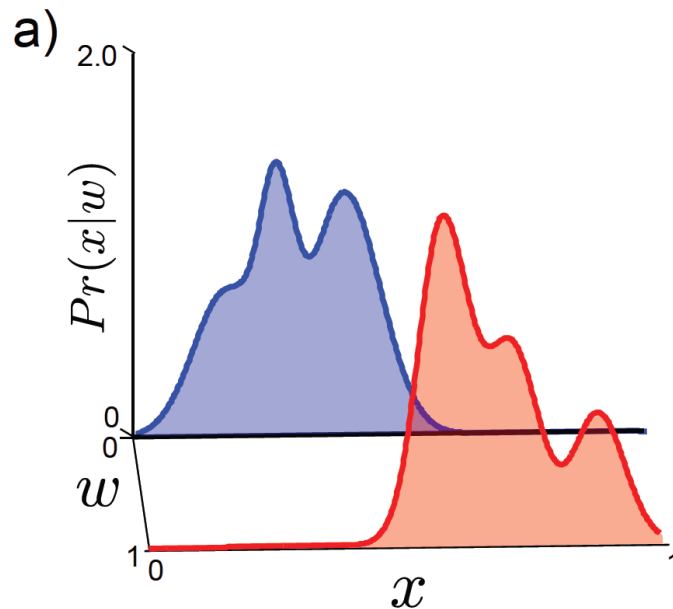
Which type should we choose?

Applications

Which type of model to use?



1. Generative methods model data – costly and many aspects of data may have no influence on world state



Which type of model to use?



2. Inference simple in discriminative models
3. Data really is generated from world – generative matches this
4. If missing data, then generative preferred
5. Generative allows imposition of prior knowledge specified by user

Structure



Computer vision models

- Three types of model

Worked example 1: Regression

Worked example 2: Classification

Which type should we choose?

Applications

Application: Skin Detection



Figure 6.7 Skin detection. For each pixel we aim to infer a label $w \in \{0, 1\}$ denoting the absence or presence of skin based on the RGB triple \mathbf{x} . Here we modeled the class conditional density functions $Pr(\mathbf{x}|w)$ as normal distributions. a) Original image. b) Log likelihood (log of data assessed under class-conditional density function) for non-skin. c) Log likelihood for skin. d) Posterior probability of belonging to skin class. e) Thresholded posterior probability $Pr(w|\mathbf{x}) > 0.5$ gives estimate of w .

Application: Background subtraction



FIAS Frankfurt Institute
for Advanced Studies



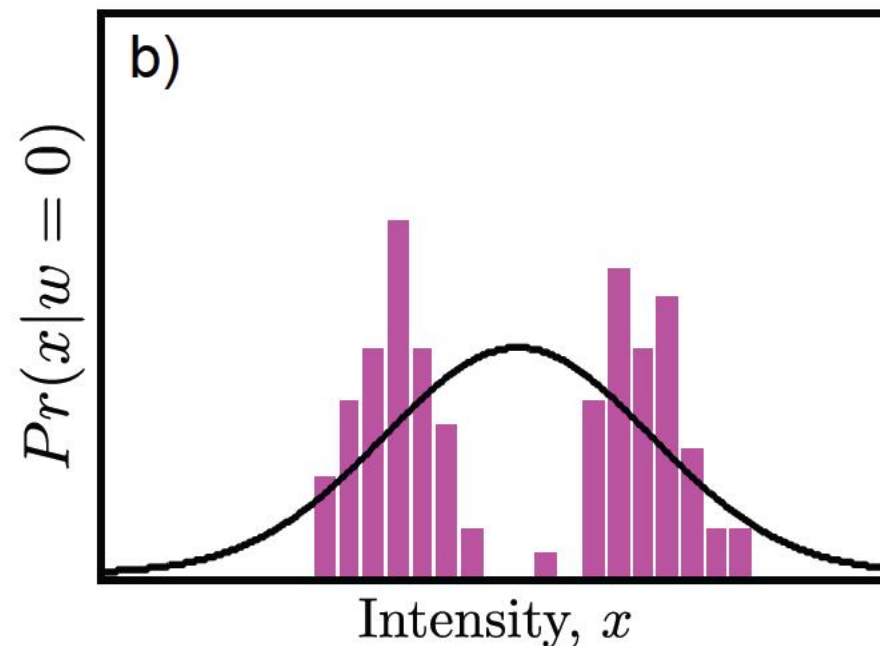
GOETHE
UNIVERSITÄT
FRANKFURT AM MAIN



Computer vision: models, learning and inference. ©2011 Simon J.D.

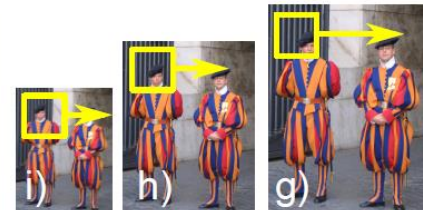
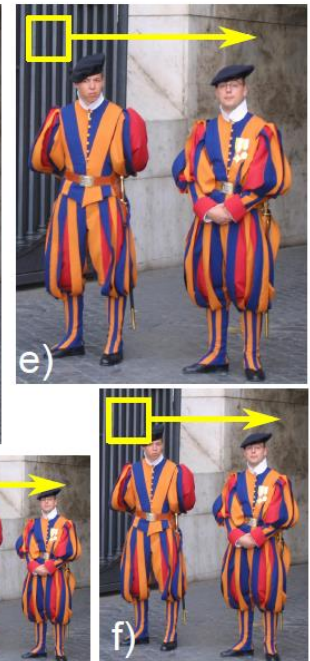
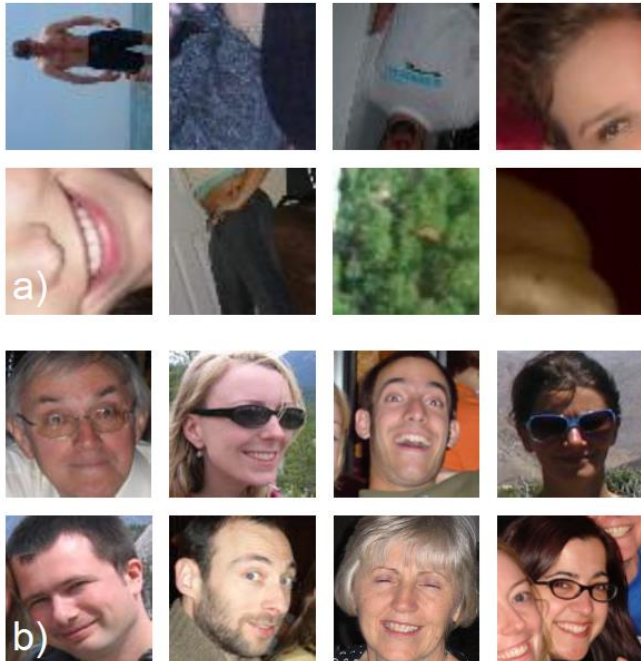
Prince

Application: Background subtraction



But consider this scene in which the foliage is blowing in the wind. A normal distribution is not good enough! Need a way to make more complex distributions

Face Detection



Type 3: $\Pr(\mathbf{x}|\mathbf{w})$ - Generative



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

- Choose an appropriate form for $\Pr(\mathbf{x})$
- Make parameters a function of \mathbf{w}
- Function takes parameters θ that define its shape

Learning algorithm: learn parameters θ from training data \mathbf{x}, \mathbf{w}

Inference algorithm: Define prior $\Pr(\mathbf{w})$ and then compute $\Pr(\mathbf{w}|\mathbf{x})$ using Bayes' rule

$$\Pr(w = 1|\mathbf{x}) = \frac{\Pr(\mathbf{x}|w = 1)\Pr(w = 1)}{\sum_{k=0}^1 \Pr(\mathbf{x}|w = k)\Pr(w = k)}$$



$$Pr(\mathbf{x}|w) = \text{Norm}_{\mathbf{x}}[\boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w]$$

Or writing in terms of class conditional density functions

$$Pr(\mathbf{x}|w = 0) = \text{Norm}_{\mathbf{x}}[\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0]$$

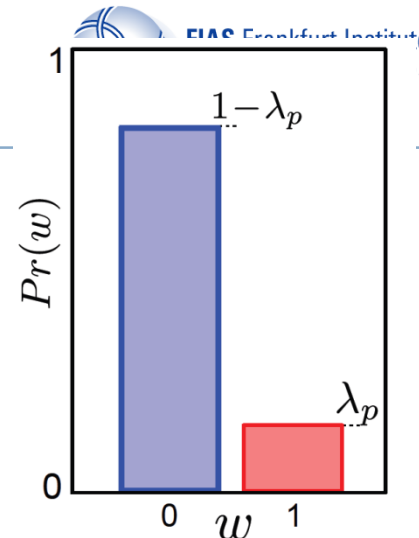
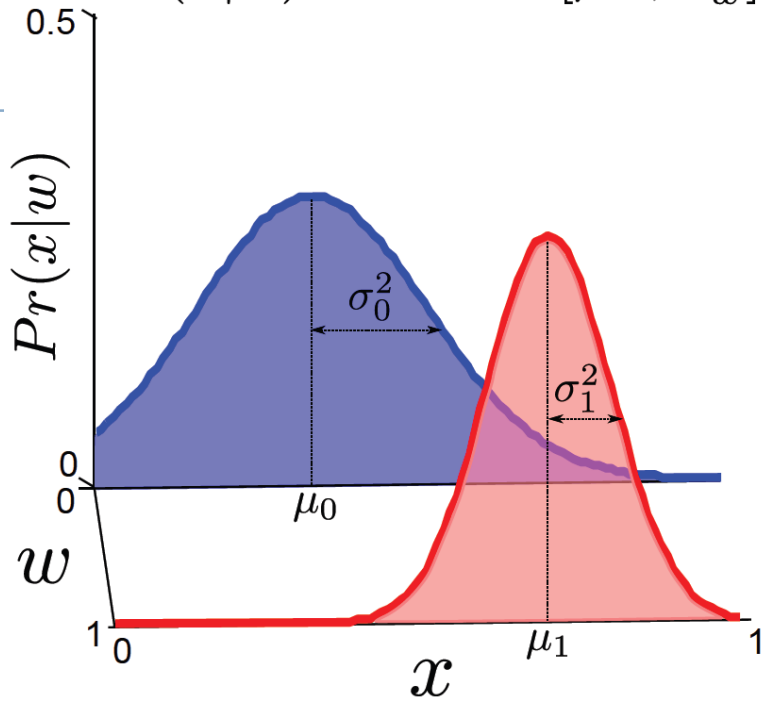
$$Pr(\mathbf{x}|w = 1) = \text{Norm}_{\mathbf{x}}[\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1]$$

Parameters $\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0$ learnt just from data S_0 where $w=0$

$$\begin{aligned}\hat{\boldsymbol{\mu}}_0, \hat{\boldsymbol{\Sigma}}_0 &= \underset{\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0}{\text{argmax}} \left[\prod_{i \in S_0} Pr(\mathbf{x}_i | \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \right] \\ &= \underset{\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0}{\text{argmax}} \left[\prod_{i \in S_0} \text{Norm}_{\mathbf{x}_i}[\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0] \right]\end{aligned}$$

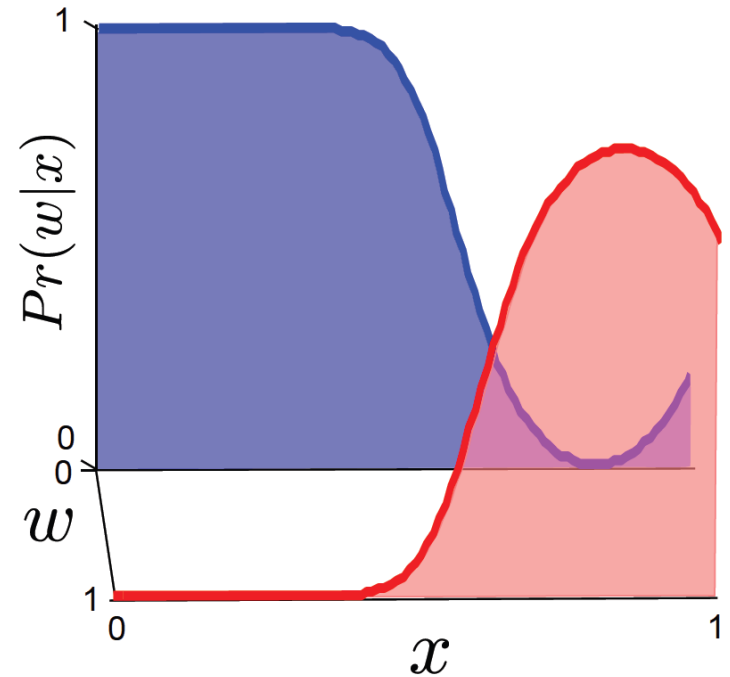
Similarly, parameters $\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1$ learnt just from data S_1 where $w=1$

$$Pr(x|w) = \text{Norm}_x[\mu_w, \sigma_w^2]$$



Inference algorithm: Define prior $Pr(\mathbf{w})$ and then compute $Pr(\mathbf{w}|\mathbf{x})$ using Bayes' rule

$$Pr(\mathbf{w}|\mathbf{x}) = \frac{Pr(\mathbf{x}|\mathbf{w})Pr(\mathbf{w})}{\int Pr(\mathbf{x}|\mathbf{w})Pr(\mathbf{w})d\mathbf{w}}$$



Experiment



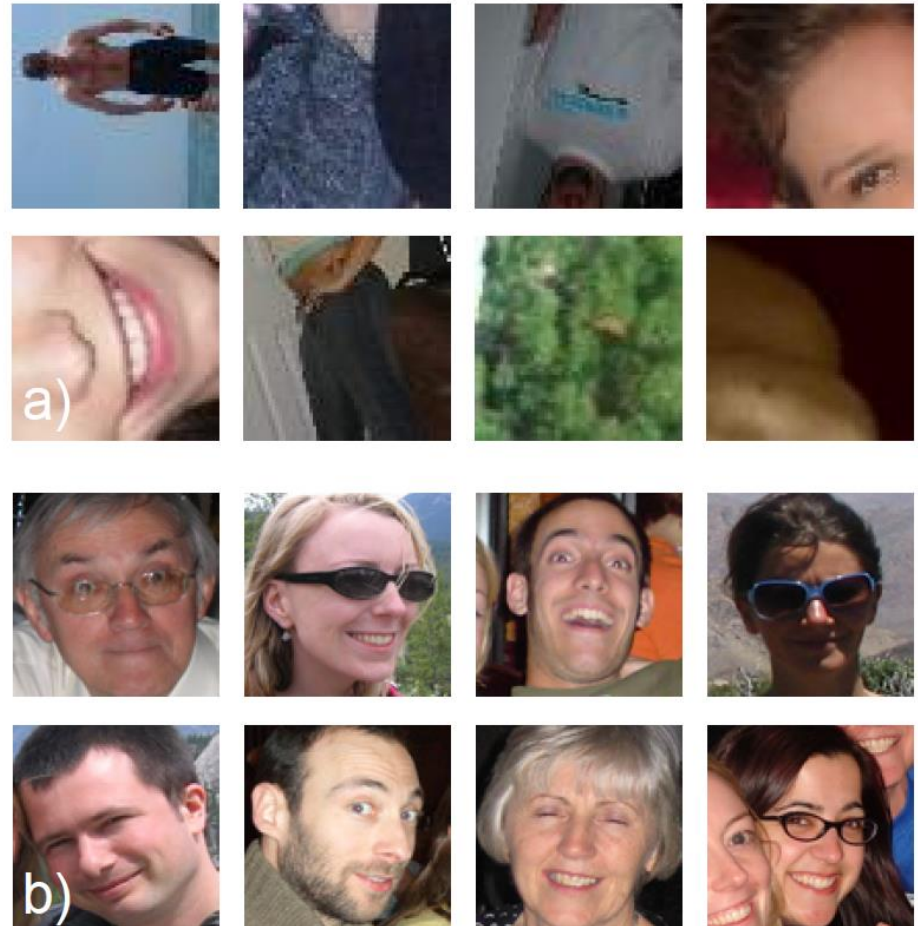
1000 non-faces

1000 faces

60x60x3 Images = 10800 x1
vectors

Equal priors $\Pr(y=1)=\Pr(y=0) =$
0.5

75% performance on test set. Not
very good!



Results (diagonal covariance)

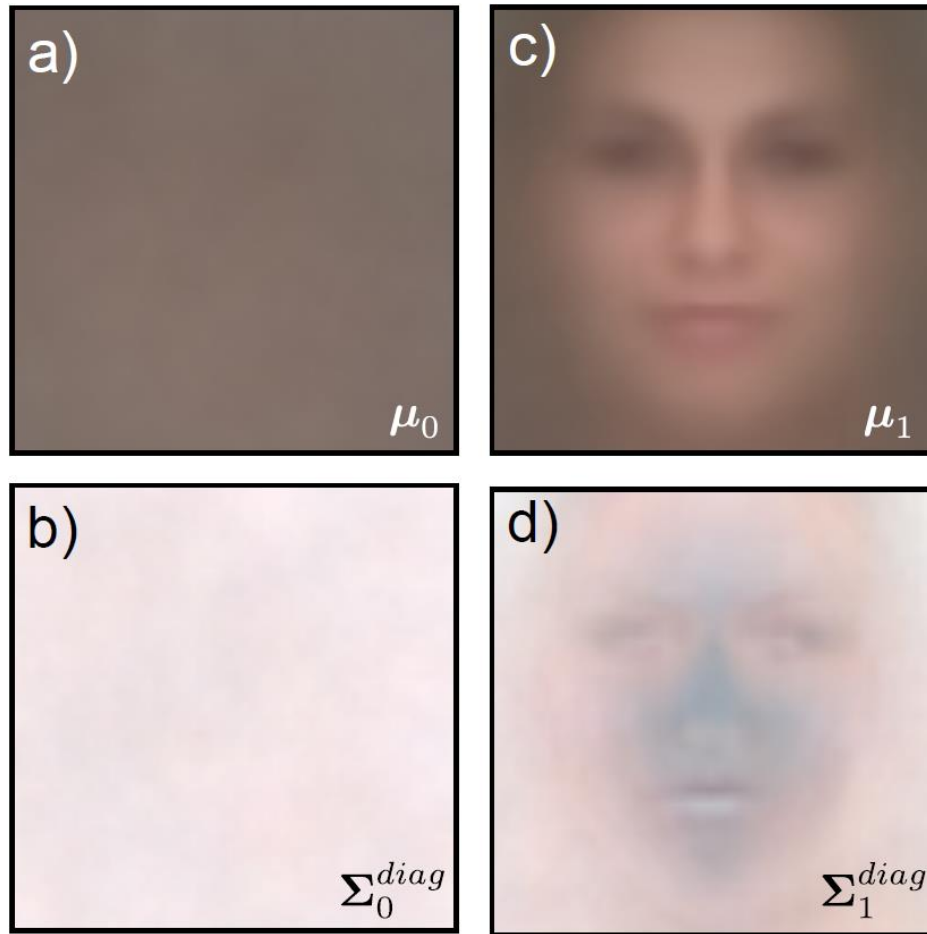


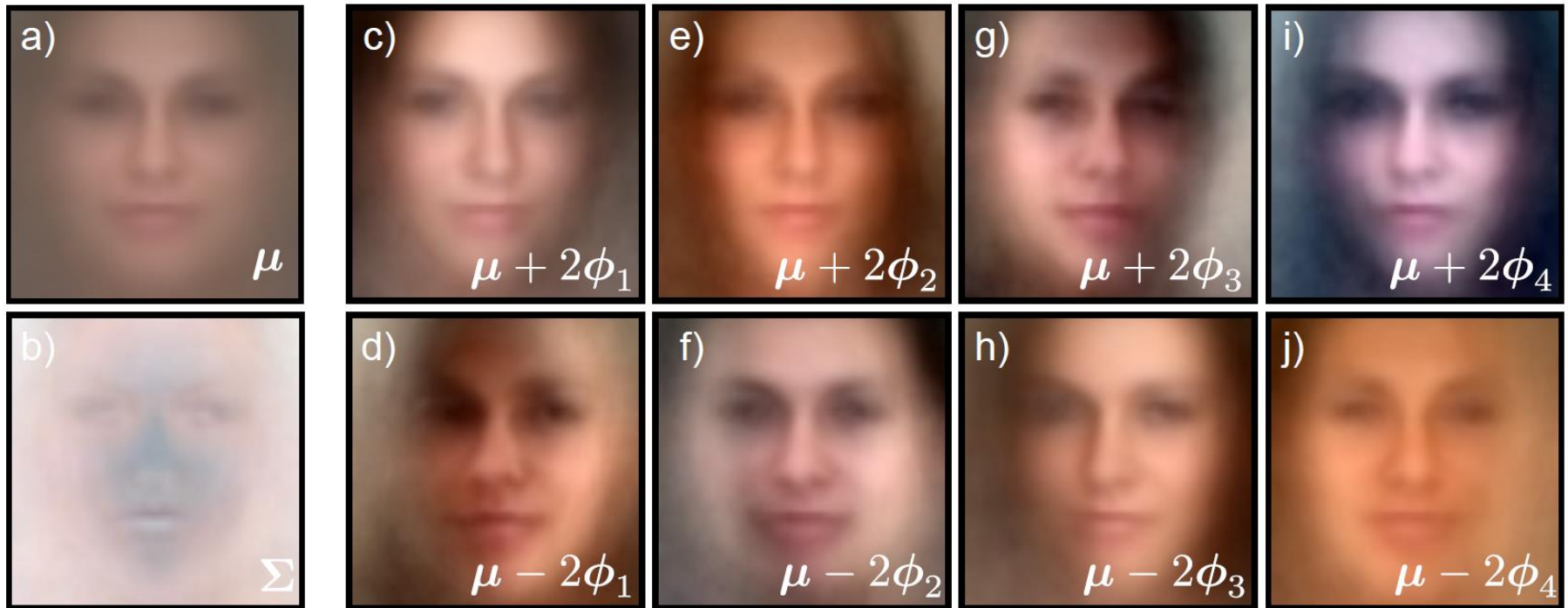
Figure 7.2 Class conditional density functions for normal model with diagonal covariance. Maximum likelihood fits based on 1000 training examples per class. a) Mean for background data μ_0 (reshaped from 10800×1 vector to 60×60 RGB image). b) Reshaped square root of diagonal covariance for background data Σ_0 . c) Mean for face data μ_1 d) Covariance for face data Σ_1 . The background model has little structure: the mean is uniform and the variance is high everywhere. The mean of the face model clearly captures class-specific information. The covariance of the face is larger at the edges of the image which usually contain hair or background.

Means of face/non-face model



Classification → 84% (9% improvement!)

Face model



Sampling from 10 parameter model

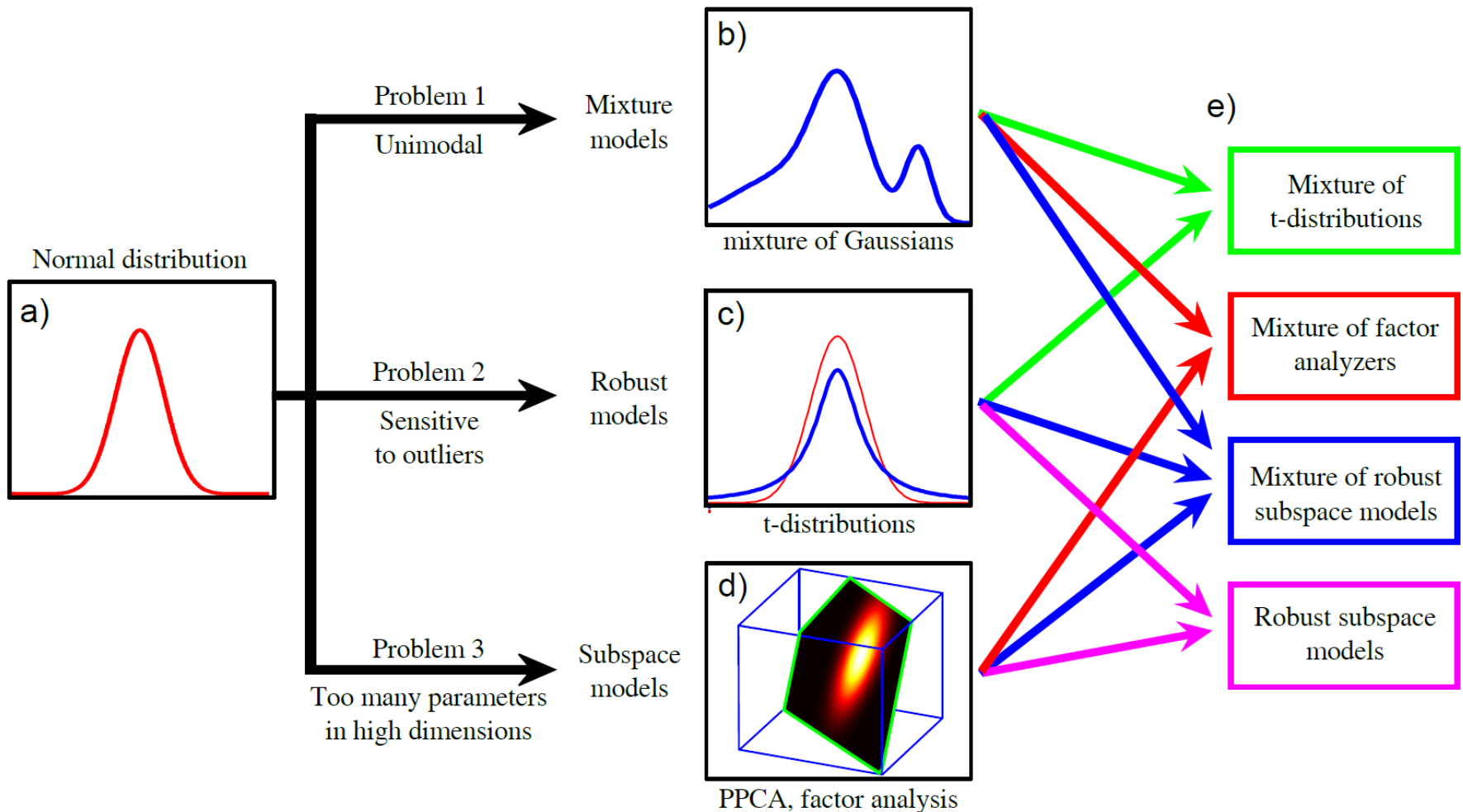
To generate:

- Choose factor loadings, \mathbf{h}_i from standard normal distribution
- Multiply by factors, Φ
- Add mean, μ
- (should add random noise component ε_i w/ diagonal cov Σ)



Computer vision: models, learning and inference. ©2011 Simon J.D.

Probability Density Models



- To do computer vision we build a model relating the image data \mathbf{x} to the world state that we wish to estimate \mathbf{w}
- Three types of model
 - Model $\Pr(\mathbf{w}|\mathbf{x})$ -- discriminative
 - Model $\Pr(\mathbf{w}|\mathbf{x})$ – generative



Backup