

Introduction to Machine Learning

A statistical perspective on supervised learning

Prof. Andreas Krause
Learning and Adaptive Systems (las.ethz.ch)

Motivation

- We have seen how we can fit prediction models (linear, non-linear) for regression and classification
- So far, these models do not have any statistical interpretation
- Often we would like to **statistically model** the data:
 - Quantify uncertainty
 - Express prior knowledge / assumptions about the data
- In the following, we will see how many of the approaches we have discussed can be interpreted as **fitting probabilistic models**
- This view will allow us to derive new methods

Minimizing the least squares error

- Assuming the data is generated iid according to

$$(\mathbf{x}_i, y_i) \sim P(\mathbf{X}, Y)$$

- The hypothesis h^* minimizing $R(h) = \mathbb{E}_{\mathbf{x},y}[(y - h(\mathbf{x}))^2]$ is given by the **conditional mean**

$$h^*(\mathbf{x}) = \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}]$$

- This (in practice unattainable) hypothesis is called the **Bayes' optimal predictor** for the squared loss

In practice we have finite data

- We know that

$$h^*(\mathbf{x}) = \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}]$$

- Thus, one strategy for estimating a predictor from training data is to estimate the conditional distribution

$$\hat{P}(Y \mid \mathbf{X})$$

and then, for test point \mathbf{x} , predict label

$$\hat{y} = \hat{\mathbb{E}}[Y \mid \mathbf{X} = \mathbf{x}] = \int \hat{P}(y \mid \mathbf{X} = \mathbf{x}) y dy$$

Estimating conditional distributions

- Common approach: Parametric estimation

- Choose a particular parametric form $\hat{P}(Y \mid \mathbf{X}, \theta)$
- Then optimize the parameters. How?

→ Maximum (conditional) Likelihood Estimation

$$\theta^* = \arg \max_{\theta} \hat{P}(y_1, \dots, y_n \mid \mathbf{x}_1, \dots, \mathbf{x}_n, \theta)$$

Ridge regression = MAP estimation

- Ridge regression can be understood as finding the Maximum A Posteriori (MAP) parameter estimate for a linear regression problem, assuming that
 - The noise $P(y|\mathbf{x}, \mathbf{w})$ is iid Gaussian and
 - The prior $P(\mathbf{w})$ on the model parameters \mathbf{w} is Gaussian

$$\arg \min_{\mathbf{w}} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_2^2 \equiv \arg \max_{\mathbf{w}} P(\mathbf{w}) \prod_i P(y_i \mid \mathbf{x}_i, \mathbf{w})$$

Regularization vs. MAP inference

- More generally, regularized estimation can often be understood as MAP inference

$$\arg \min_{\mathbf{w}} \sum_{i=1}^n \ell(\mathbf{w}^T \mathbf{x}_i; \mathbf{x}_i, y_i) + C(\mathbf{w}) = \arg \max_{\mathbf{w}} \prod_i P(y_i \mid \mathbf{x}_i, \mathbf{w}) \underline{P(\mathbf{w})}$$
$$= \arg \max_{\mathbf{w}} P(\mathbf{w} \mid D)$$

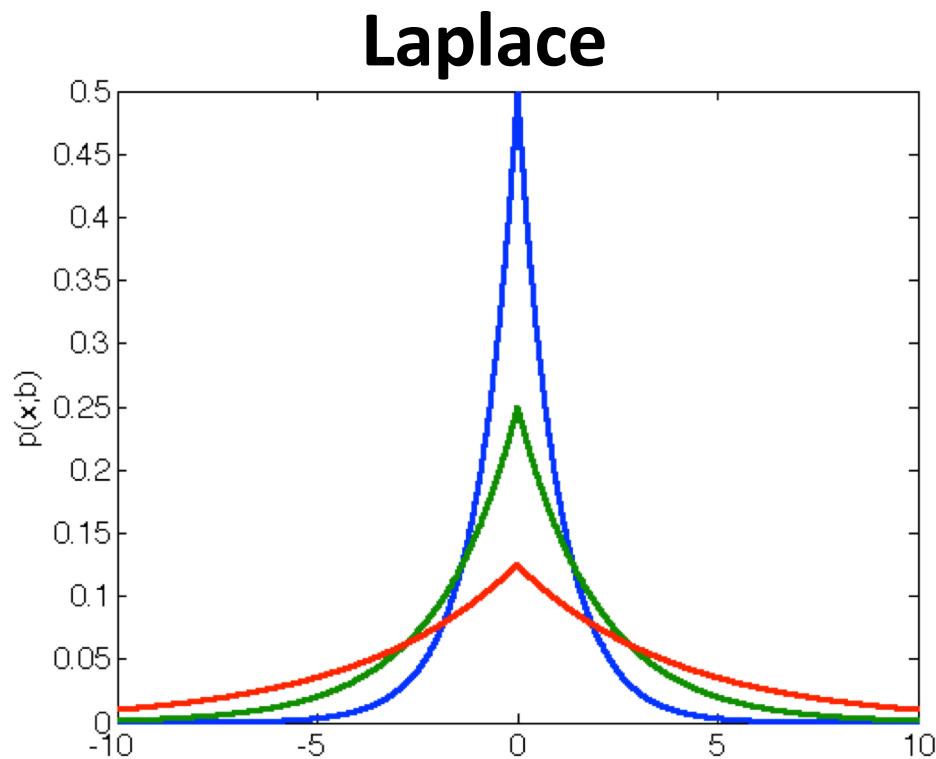
where $C(\mathbf{w}) = -\log P(\mathbf{w})$

and $\ell(\mathbf{w}^T \mathbf{x}_i; \mathbf{x}_i, y_i) = -\log P(y_i \mid \mathbf{x}_i, \mathbf{w})$

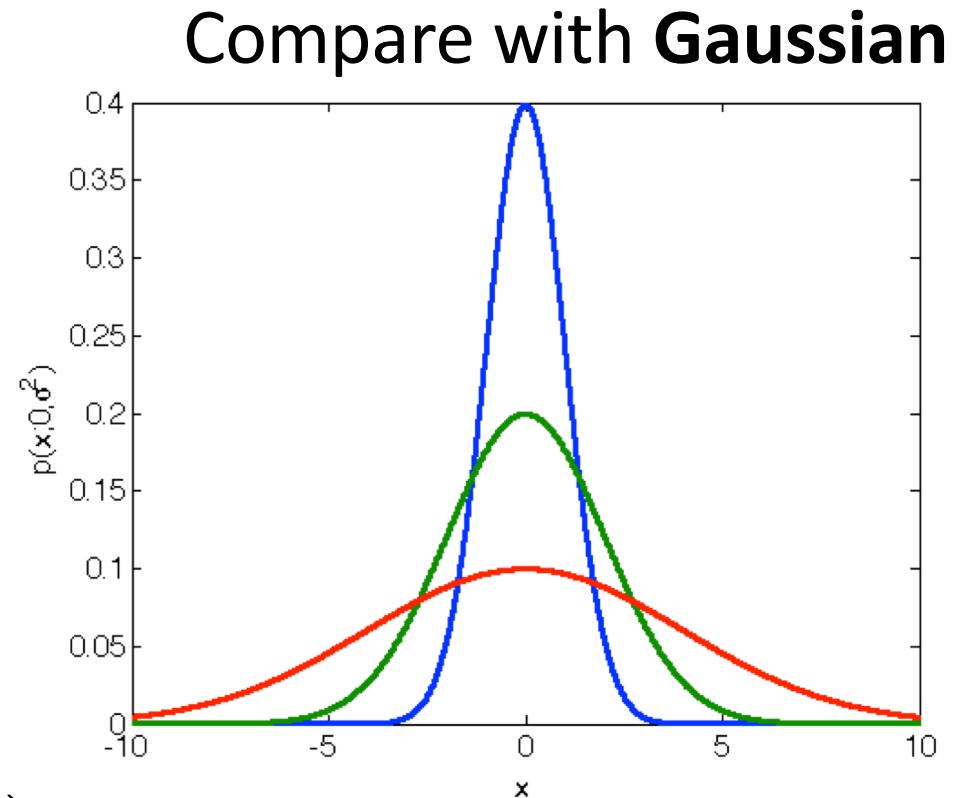
- This perspective allows changing priors (=regularizers) and likelihoods (=loss functions)

Example: l1-regularization

- Is there a prior that corresponds to l1-regularization?
- **Answer:** The [Laplace prior](#)



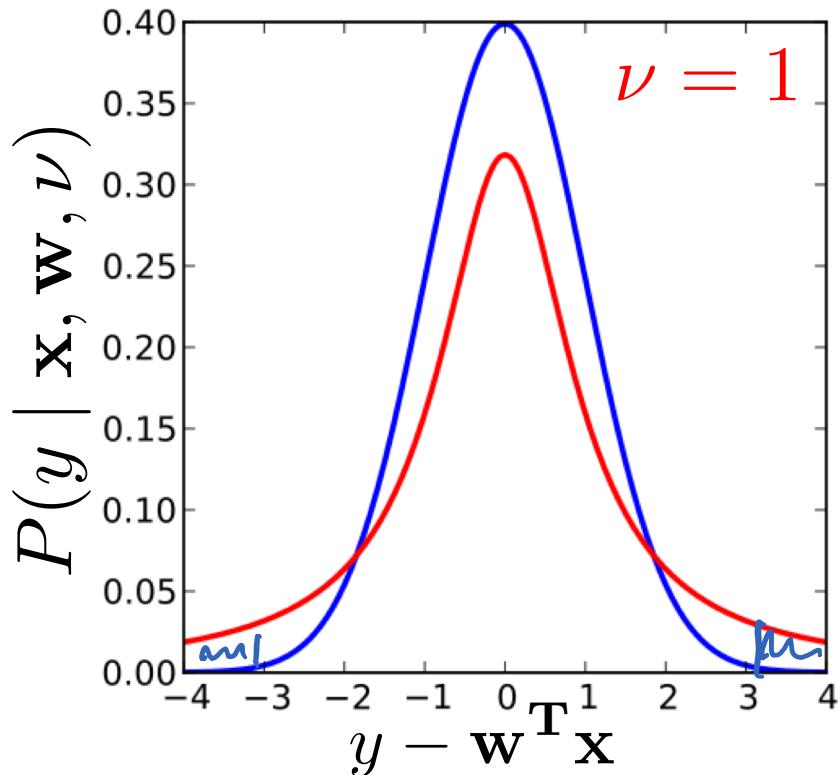
$$p(x; \mu, b) = \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right)$$



Example: student-t likelihood

- Can introduce **robustness** by changing the likelihood (=loss) function
- **Example:** (non-standardized) Student's-t likelihood

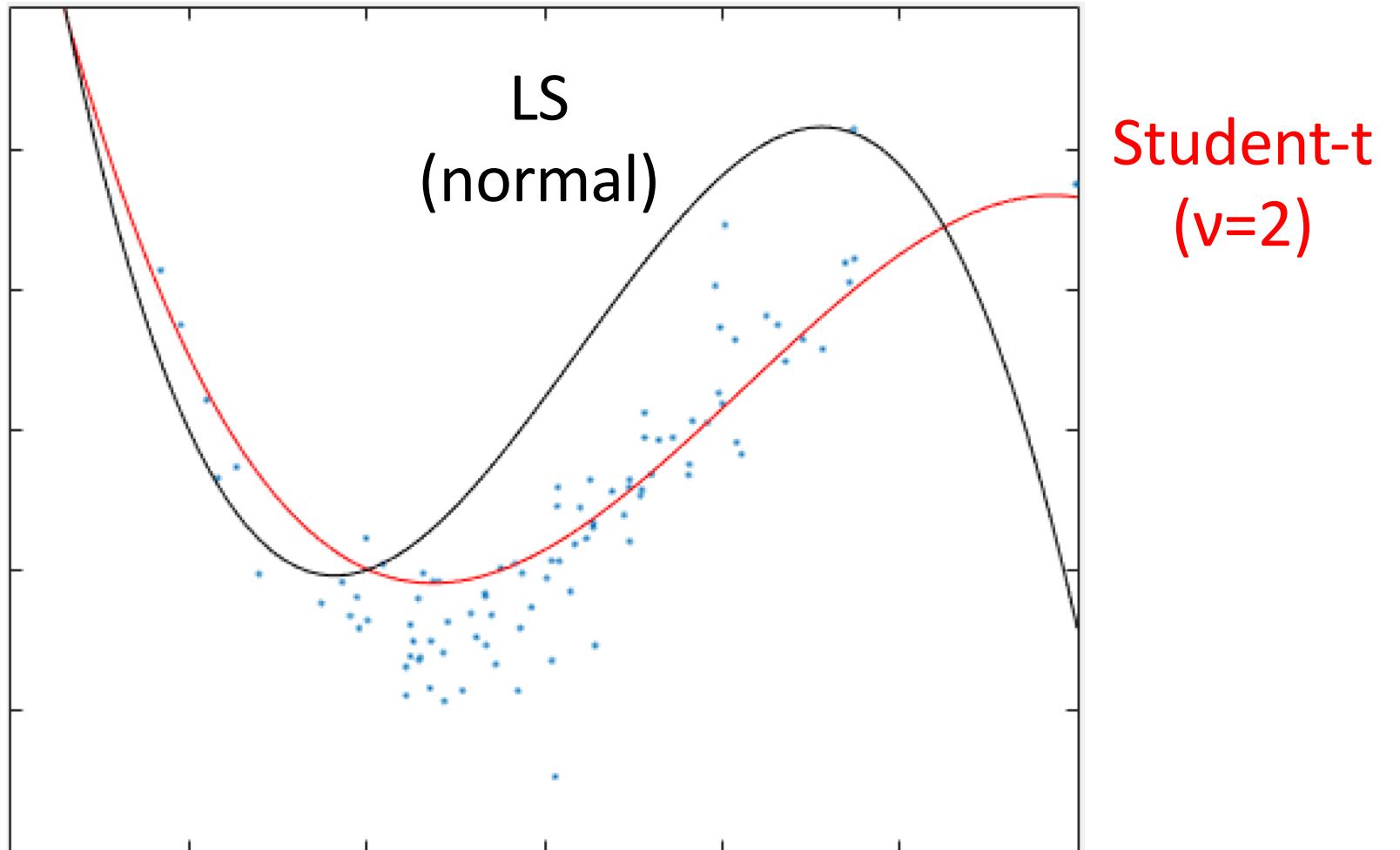
$$P(y \mid \mathbf{x}, \mathbf{w}, \nu, \sigma^2) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\pi\nu\sigma^2}\Gamma(\frac{\nu}{2})} \left(1 + \frac{(y - \mathbf{w}^T \mathbf{x})^2}{\nu\sigma^2}\right)^{-\frac{\nu+1}{2}}$$



For Gaussian: $P(|y - \mathbf{w}^T \mathbf{x}| > t \cdot \sigma) = O(e^{-t})$

For student-t: $\sim e^{-c} \cdot \dots = O(t^{-\alpha})$

Example fits



Statistical models for classification

- So far, we have focused on *regression*
- Are there natural statistical models for *classification*?

Risk in classification

- In classification, risk is

$$R(h) = \mathbb{E}_{\mathbf{X}, Y} [[Y \neq h(\mathbf{X})]]$$

$\left\{ \begin{array}{ll} 1 & \text{if } Y \neq h(x) \\ 0 & \text{otw.} \end{array} \right.$

- Suppose (unrealistically) we knew $P(X, Y)$
- Which h minimizes the risk then?

$$h^*(x) = \underset{\hat{y}}{\operatorname{argmin}} \mathbb{E}_Y [[Y \neq \hat{y}] | X=x] = (\star)$$

$\ell(\hat{y})$

$$\ell(\hat{y}) = \sum_{y=1}^C P(Y=y | X=x) [y \neq \hat{y}] = \sum_{y: y \neq \hat{y}} P(Y=y | X=x) = 1 - P(Y=\hat{y} | X=x)$$

$$(\star) = \underset{\hat{y}}{\operatorname{argmax}} P(Y=\hat{y} | X=x)$$

Bayes' optimal classifier

- Assuming the data is generated iid according to

$$(\mathbf{x}_i, y_i) \sim P(\mathbf{X}, Y)$$

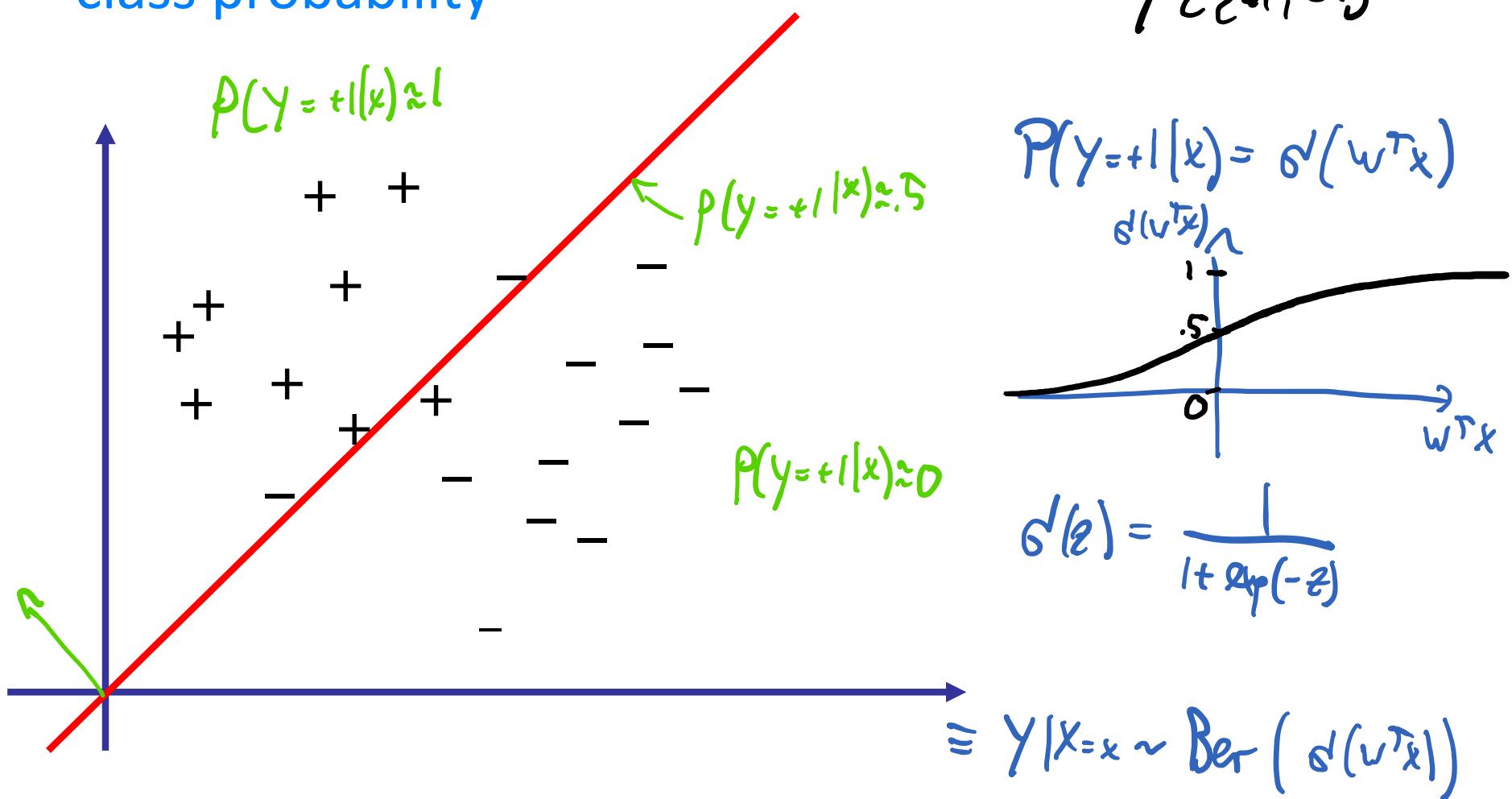
- The hypothesis h^* minimizing $R(h) = \mathbb{E}_{\mathbf{X}, Y}[[Y \neq h(\mathbf{X})]]$ is given by the **most probable class**

$$h^*(x) = \arg \max_y P(Y = y \mid \mathbf{X} = \mathbf{x})$$

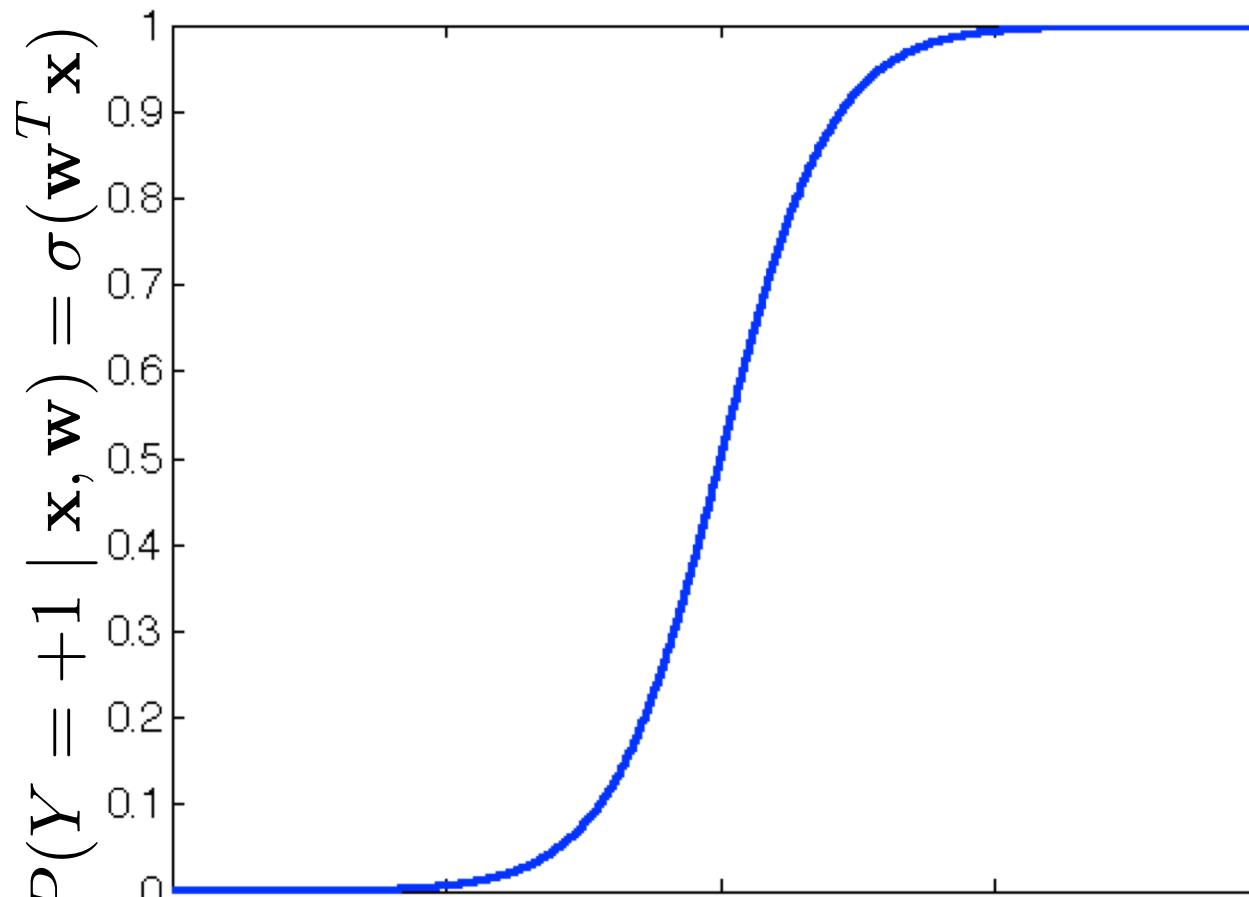
- This (in practice unattainable) hypothesis is called the **Bayes' optimal predictor** for the squared loss
- Thus, natural approach is again to estimate $P(Y \mid X)$

Logistic regression

- Idea: Use (generalized) linear model for the class probability



Link function for logistic regression



$$\begin{aligned} P(Y=+1 | \mathbf{x}, \mathbf{w}) &= \sigma(\mathbf{w}^T \mathbf{x}) \\ &\approx 1 - P(Y=-1 | \mathbf{x}) \\ &= 1 - \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})} \\ &= \frac{\exp(-\mathbf{w}^T \mathbf{x})}{1 + \exp(-\mathbf{w}^T \mathbf{x})} \\ &= \frac{1}{1 + \exp(\mathbf{w}^T \mathbf{x})} \\ &\Rightarrow P(Y=y | \mathbf{x}) \\ &\approx \frac{1}{1 + \exp(-y \cdot \mathbf{w}^T \mathbf{x})} \end{aligned}$$

- Link function

$$\sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

Logistic regression

- Logistic regression (a classification method) replaces the assumption of Gaussian noise (squared loss) by iid Bernoulli noise:

$$P(y \mid \mathbf{x}, \mathbf{w}) = \text{Ber}(y; \sigma(\mathbf{w}^T \mathbf{x}))$$

vs. $P(y_i \mid \mathbf{x}_i, \mathbf{w}) = \mathcal{N}(y_i; \mathbf{w}^T \mathbf{x}_i, \sigma^2)$ for LS regression

- How can we estimate the parameters \mathbf{w} ?
- Maximum Likelihood Estimation / MAP estimation

MLE for logistic regression

$$\hat{w} \in \arg\max_w P(D|w) \stackrel{iid}{=} \arg\max_w \prod_{i=1}^n P(y_i|x_i, w)$$

$$= \arg\min_w - \sum_{i=1}^n \log \underline{P(y_i|x_i, w)}$$

$$\begin{aligned} P(y_i|x_i, w) &= \frac{P(g_i; x_i|w)}{P(x_i|w)} \\ &= c_i P(y_i; x_i|w) \\ &\quad \text{if } P(x_i|w) = P(x_i) \\ &\quad \text{i.e. } x_i \text{ independent} \end{aligned}$$

$$-\log P(y|x, w) = -\log \frac{1}{1 + \exp(-y w^T x)} = \log (1 + \exp(-y w^T x))$$

$$(R) = \arg\min_w \sum_{i=1}^n \boxed{\log (1 + \exp(-y_i w^T x_i))}$$

$\ell_{\text{logistic}}(w; x_i, y_i)$

$R(w)$

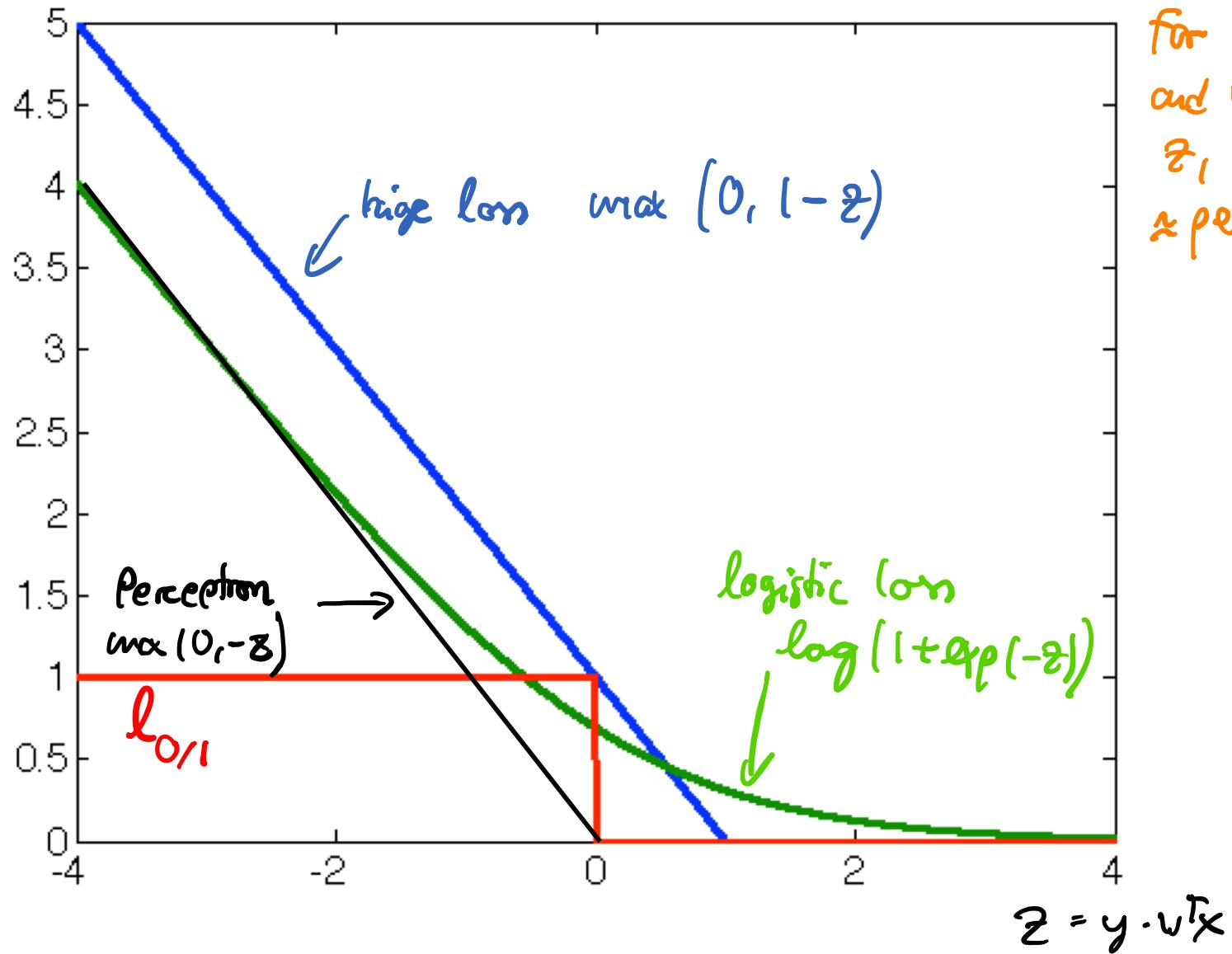
MLE for logistic regression

- Negative log likelihood (=objective) function given by

$$\hat{R}(\mathbf{w}) = \sum_{i=1}^n \log\left(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i)\right)$$

- The **logistic loss** is **convex**
→ Can use convex optimization techniques (e.g., SGD)

Logistic loss vs other losses



Gradient for logistic regression

- Loss for data point (\mathbf{x}, y)

$$\ell(\mathbf{w}) = \log\left(1 + \exp(-y\mathbf{w}^T \mathbf{x})\right)$$

$$\begin{aligned}\nabla_{\mathbf{w}} \ell(\mathbf{w}) &= \frac{1}{1 + \exp(-y\mathbf{w}^T \mathbf{x})} \cdot \exp(-y\mathbf{w}^T \mathbf{x}) \cdot (-y \cdot \mathbf{x}) \\ &= \frac{\exp(-y\mathbf{w}^T \mathbf{x})}{1 + \exp(-y\mathbf{w}^T \mathbf{x})} \cdot (-y \cdot \mathbf{x}) \\ &= \underbrace{\frac{1}{1 + \exp(y \cdot \mathbf{w}^T \mathbf{x})}}_{P(y \neq y | \mathbf{x})} \cdot (-y \cdot \mathbf{x})\end{aligned}$$

SGD for logistic regression

- Initialize \mathbf{w}
- For $t = 1, 2, \dots$
 - Pick data point (\mathbf{x}, y) uniformly at random from data D
 - Compute probability of misclassification with current model

$$\hat{P}(Y = -y \mid \mathbf{w}, \mathbf{x}) = \frac{1}{1 + \exp(y\mathbf{w}^T \mathbf{x})}$$

- Take gradient step

$$\mathbf{w} \leftarrow \mathbf{w} + \eta_t y \mathbf{x} \hat{P}(Y = -y \mid \mathbf{w}, \mathbf{x})$$

Logistic regression and regularization

- Similar to SVMs and linear regression, want to use **regularizer** to control model complexity
- Thus, instead of solving MLE

$$\min_{\mathbf{w}} \sum_{i=1}^n \log\left(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i)\right)$$

estimate MAP/solve regularized problem

- **L2 (Gaussian prior):**

$$\min_{\mathbf{w}} \sum_{i=1}^n \log\left(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i)\right) + \underline{\lambda \|\mathbf{w}\|_2^2}$$

- **L1 (Laplace):**

$$\min_{\mathbf{w}} \sum_{i=1}^n \log\left(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i)\right) + \underline{\lambda \|\mathbf{w}\|_1}$$

SGD for L2-regularized logistic regression

- Initialize \mathbf{w}
- For $t = 1, 2, \dots$
 - Pick data point (\mathbf{x}, y) uniformly at random from data D
 - Compute probability of misclassification with current model

$$\hat{P}(Y = -y \mid \mathbf{w}, \mathbf{x}) = \frac{1}{1 + \exp(y\mathbf{w}^T \mathbf{x})}$$

- Take gradient step

$$\mathbf{w} \leftarrow \mathbf{w} (1 - 2\lambda\eta_t) + \eta_t y \mathbf{x} \hat{P}(Y = -y \mid \mathbf{w}, \mathbf{x})$$


Regularized logistic regression

- Learning:

- Find optimal weights by minimizing logistic loss + regularizer

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{i=1}^n \log \left(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i) \right) + \underline{\lambda \|\mathbf{w}\|_2^2}$$
$$= \arg \max_{\mathbf{w}} P(\mathbf{w} \mid \mathbf{x}_1, \dots, \mathbf{x}_n, y_1, \dots, y_n)$$

- Classification:

- Use conditional distribution

$$P(y \mid \mathbf{x}, \hat{\mathbf{w}}) = \frac{1}{1 + \exp(-y \hat{\mathbf{w}}^T \mathbf{x})}$$

- E.g., predict more likely class label

$$(\#) \quad \hat{y} = \text{Sign}(\hat{\mathbf{w}}^T \mathbf{x})$$

(*) compute $P(\hat{y} \mid \mathbf{x}, \hat{\mathbf{w}})$

\hat{y}

$= \text{compute } \frac{1}{1 + \exp(-\hat{y} \hat{\mathbf{w}}^T \mathbf{x})}$

$= \underset{y}{\text{argmax}} \exp(-\hat{y} \hat{\mathbf{w}}^T \mathbf{x})$

$= \underset{\hat{y}}{\text{argmin}} -\hat{y} \hat{\mathbf{w}}^T \mathbf{x}$

$= \underset{y \in \{-1, 1\}}{\text{argmax}} \hat{y} \hat{\mathbf{w}}^T \mathbf{x} = \text{Sgn}(\hat{\mathbf{w}}^T \mathbf{x})$

Logistic regression demo

More remarks on logistic regression

- Can kernelize ([kernelized logistic regression](#))
- Can apply logistic loss function to neural networks,
in order to have them output probabilities
- Natural multi-class variants
- ...

Kernelized logistic regression

- Learning:

- Find optimal weights by minimizing logistic loss + regularizer

$$\hat{\alpha} = \arg \min_{\alpha} \sum_{i=1}^n \log \left(1 + \exp(-y_i \alpha^T \mathbf{K}_i) \right) + \lambda \alpha^T \mathbf{K} \alpha$$
$$\mathbf{K} = \begin{pmatrix} k_1 & | & \dots & | & k_n \end{pmatrix}$$

- Classification:

- Use conditional distribution

$$\hat{P}(y \mid \mathbf{x}, \hat{\alpha}) = \frac{1}{1 + \exp \left(-y \sum_{j=1}^n \alpha_j k(\mathbf{x}_j, \mathbf{x}) \right)}$$

- E.g., predict more likely class label

$w^T x$ assuming
 $w = \sum_i \alpha_i k_i$

Multi-class logistic regression

- Can extend logistic regression to multi-class setting
- Maintain one weight vector per class and model

$$P(Y = i \mid \mathbf{x}, \mathbf{w}_1, \dots, \mathbf{w}_c) = \frac{\exp(\mathbf{w}_i^T \mathbf{x})}{\sum_{j=1}^c \exp(\mathbf{w}_j^T \mathbf{x})}$$

\hat{p}_i

$\sum_{i=1}^c \hat{p}_i = 1$

$$\frac{\exp(\mathbf{w}_i^T \mathbf{x} + \gamma)}{\sum_j \exp(\mathbf{w}_j^T \mathbf{x} + \gamma)} = \frac{\exp(\mathbf{v}_i^T \mathbf{x}) \cdot e^\gamma}{\sum_j \exp(\mathbf{w}_j^T \mathbf{x}) \cdot e^\gamma}$$

\Rightarrow w.l.o.g., can pick γ to be $-\mathbf{w}_c^T \mathbf{x}$

$\hookrightarrow \mathbf{w}_i \mapsto \mathbf{v}_i - \mathbf{w}_c$, and $\mathbf{v}_c \mapsto 0$

Multi-class logistic regression

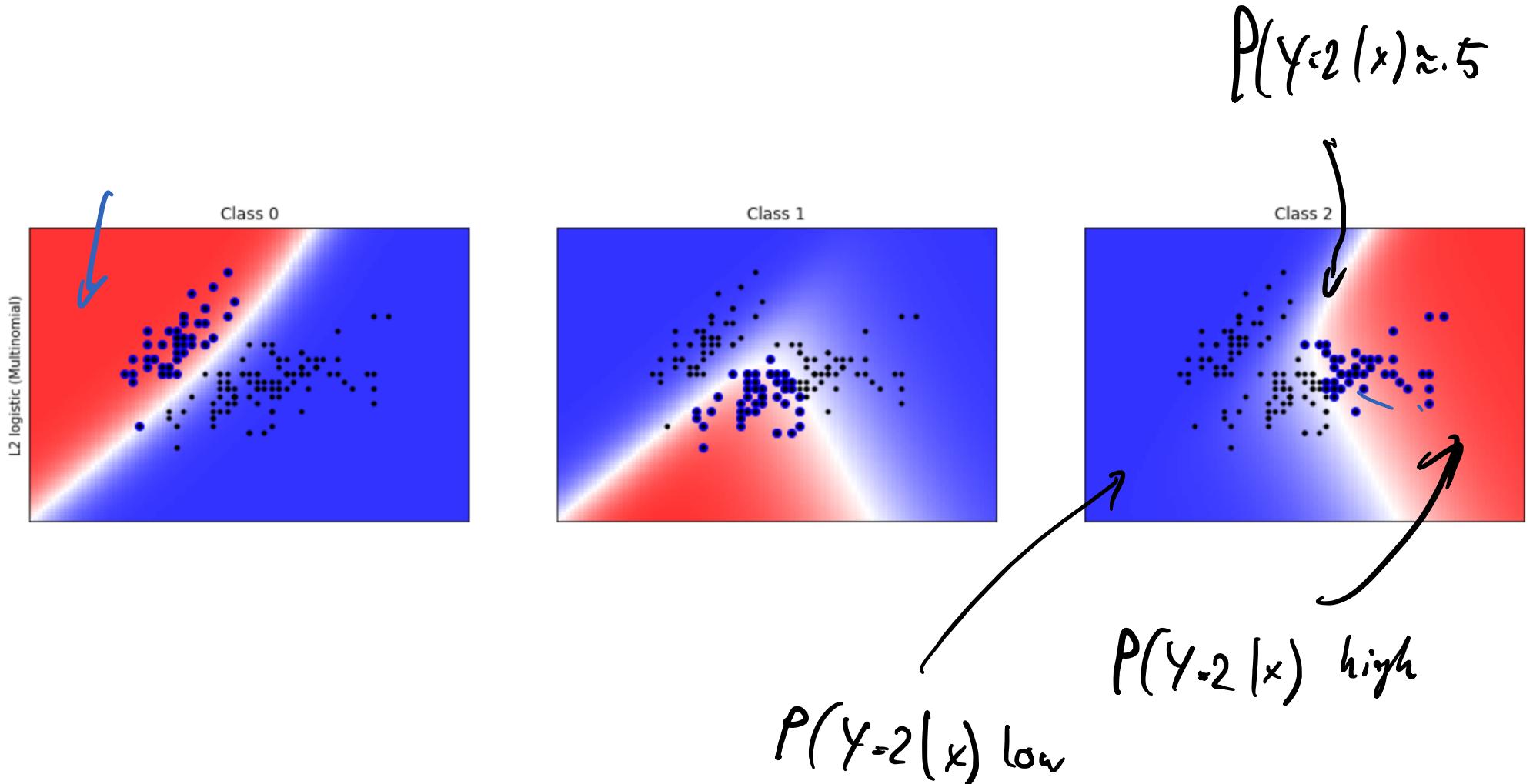
- Maintain one weight vector per class and model

$$P(Y = i \mid \mathbf{x}, \mathbf{w}_1, \dots, \mathbf{w}_c) = \frac{\exp(\mathbf{w}_i^T \mathbf{x})}{\sum_{j=1}^c \exp(\mathbf{w}_j^T \mathbf{x})}$$

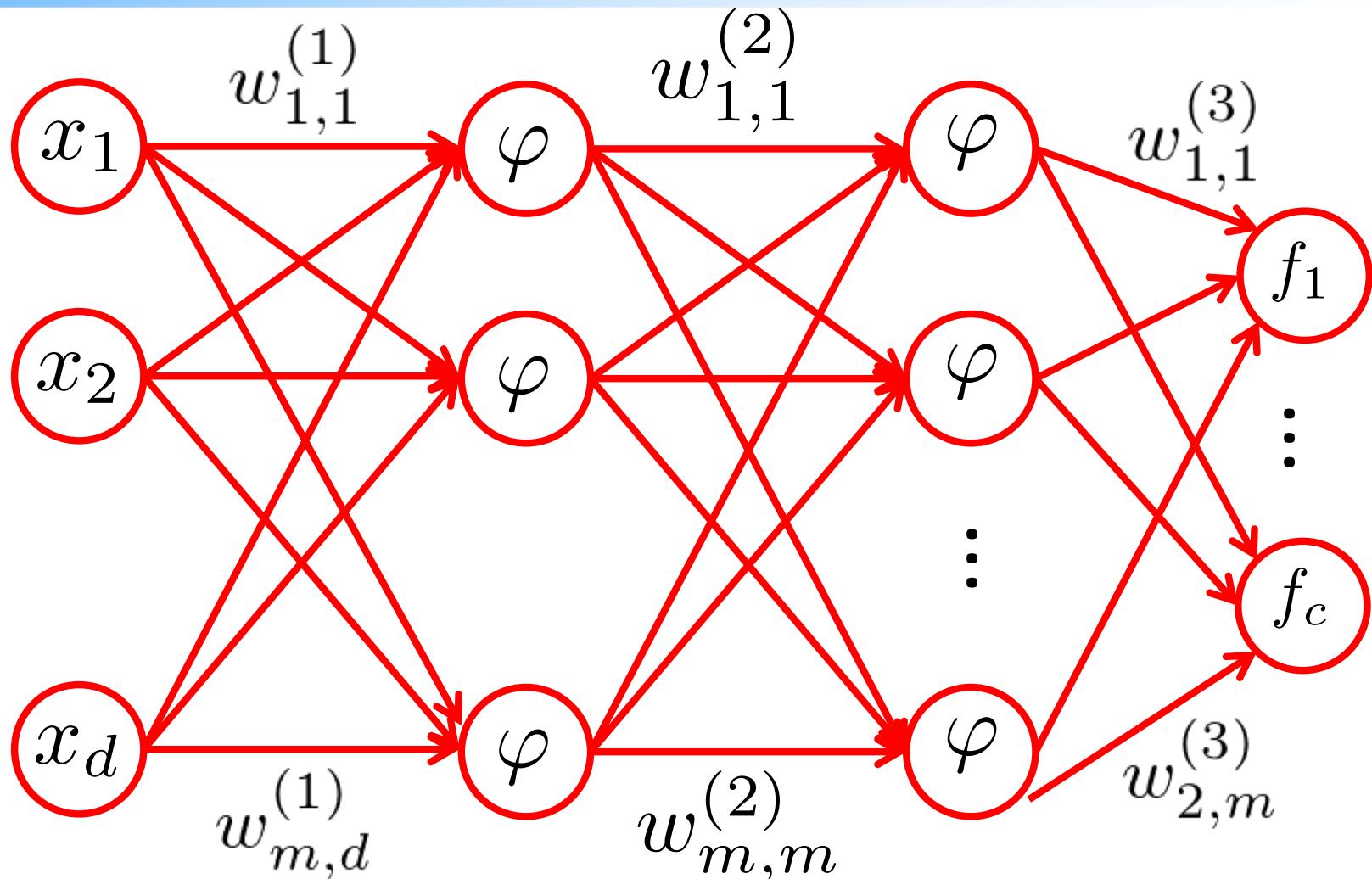
- Not unique – can force uniqueness by setting $\mathbf{w}_c = 0$ (this recovers logistic regression as special case)
- Corresponding loss function (**cross-entropy loss**):

$$\ell(y; \mathbf{x}, \mathbf{w}_1, \dots, \mathbf{w}_c) = -\log P(Y = y \mid \mathbf{x}, \mathbf{w}_1, \dots, \mathbf{w}_c)$$

Illustration



Training neural nets for multi-class



Loss: $\ell(Y = i; f_1, \dots, f_c) = -\log \frac{\exp(f_i)}{\sum_{j=1}^c \exp(f_j)}$

SVM vs. Logistic regression

<i>Method</i>	<i>SVM / Perceptron</i>	<i>Logistic regression</i>
Advantages	Sometimes higher classification accuracy; Sparse sol's	Can obtain class probabilities
Disadvantages	Can't (easily) get class probabilities	Dense solutions

Outlook: Bayesian learning

“Optimization” based learning (MAP, MLE, ...):

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} P(\mathbf{w} \mid D) \quad P(y \mid \mathbf{x}, \hat{\mathbf{w}})$$

Ignores uncertainty in model

Optimization typically efficient

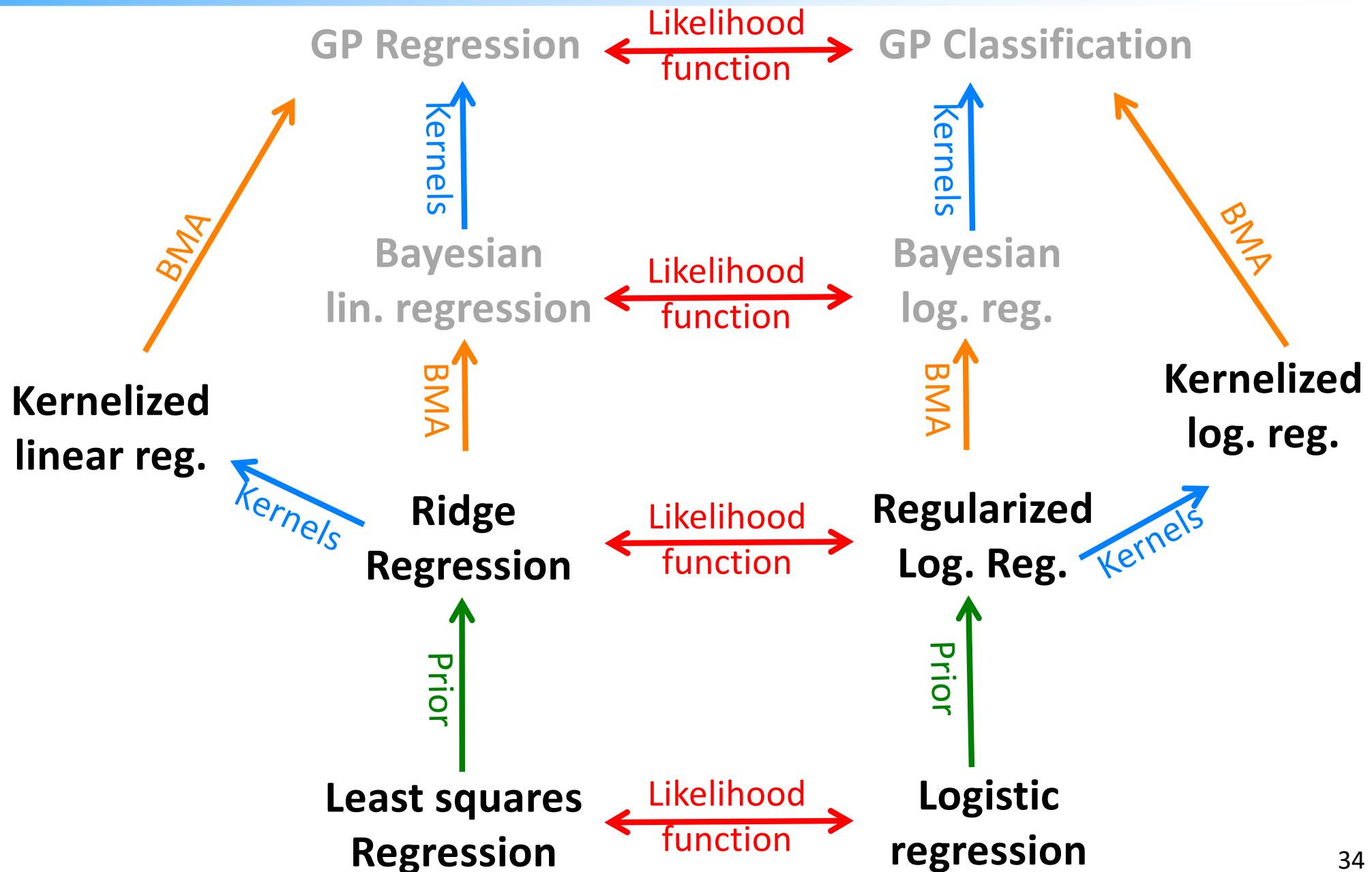
“Integration” based learning / Bayesian model averaging:

$$P(y \mid \mathbf{x}, D) = \int P(y \mid \mathbf{x}, \mathbf{w}) \underbrace{P(\mathbf{w} \mid D)}_{}$$

Quantifies uncertainty in model

integration typically intractable

Probabilistic modeling big picture so far



Representation/
features

Linear hypotheses; nonlinear hypotheses with
nonlinear feature transforms, kernels, learn nonlinear
features via neural nets

Probabilistic /
Optimization
Model:

Likelihood * Prior
Loss-function + Regularization

Squared loss = Gaussian lik., 0/1,
Perceptron, Hinge, cost sensitive,
multi-class hinge, reconstruction
error, logistic loss=Bernoulli lik.,
cross-entropy loss=Categorical lik.

L^2 norm (=Gaussian prior),
 L^1 norm (=Laplace prior),
early stopping, dropout

Method:

Exact solution, Gradient Descent, (mini-batch) SGD,
Reductions, Lloyd's heuristic, Bayesian model averaging

Evaluation
metric:

Mean squared error, Accuracy, F1 score, AUC,
Confusion matrices, compression performance,
log-likelihood on validation set

Model selection:

K-fold Cross-Validation, Monte Carlo CV,
Bayesian model selection