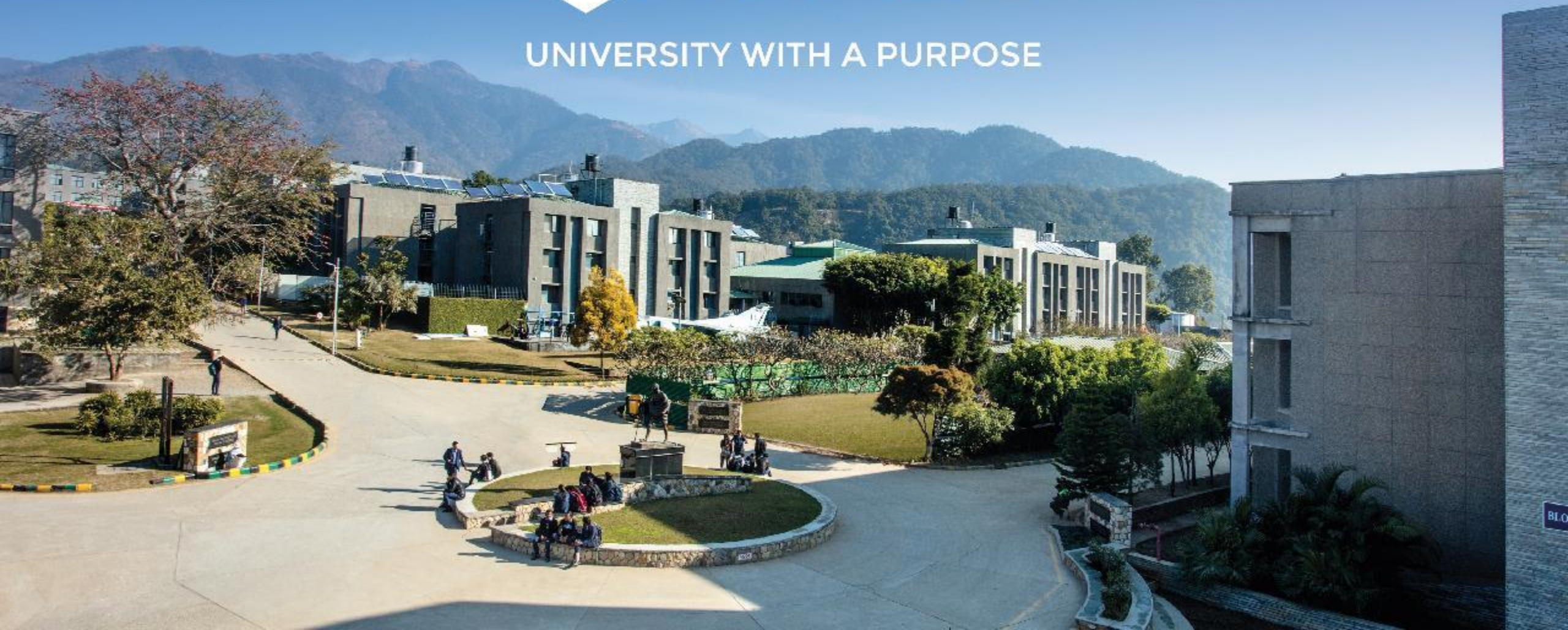
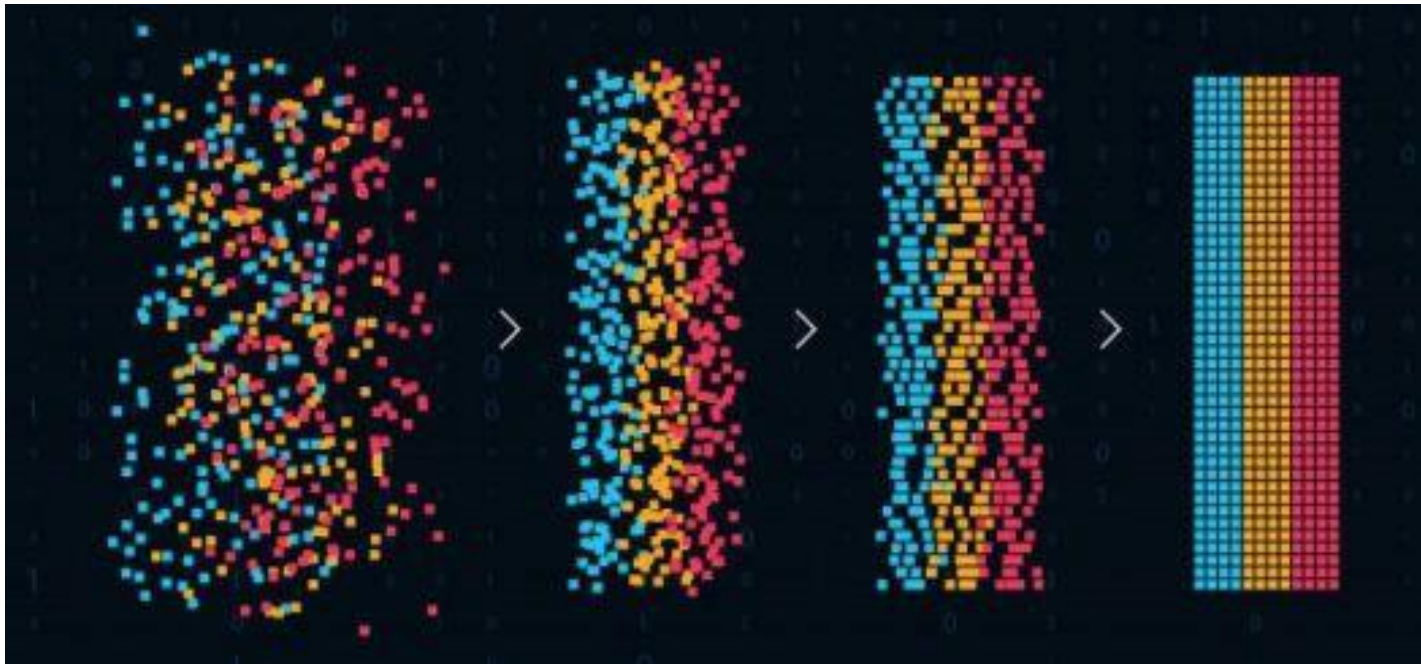




UNIVERSITY WITH A PURPOSE



# Pattern and Anomaly Detection



Source: Edureka

**B. Tech., CSE + AI/ML**

Dr Gopal Singh Phartiyal

16/09/2021



# Sequential Learning

- **Context:** batch techniques (ex. ML) takes all the training data in one go.
- This increases the computational cost and also dependency on presence of whole data at once.
- In sequential learning: Data items considered one at a time (a.k.a. on-line learning); Ex.: **stochastic (sequential) gradient descent:**

$$\begin{aligned} \mathbf{w}^{(\tau+1)} &= \mathbf{w}^{(\tau)} - \eta \nabla E_n \\ &= \mathbf{w}^{(\tau)} + \eta (t_n - \mathbf{w}^{(\tau)T} \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n). \end{aligned}$$

$\tau$  denotes the iteration number, and  $\eta$  is a learning rate parameter

- For sum-of-squares error  $\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \eta (t_n - \mathbf{w}^{(\tau)T} \phi_n) \phi_n$
- This is known as the *least-mean-squares (LMS) algorithm*.

# Regularized Least Squares

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- **Remember:** Weights being very large
- The idea of adding regularization term

- In general

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

Data term + Regularization term

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

Weight Decay regularizer

*In statistics: parameter shrinkage method*

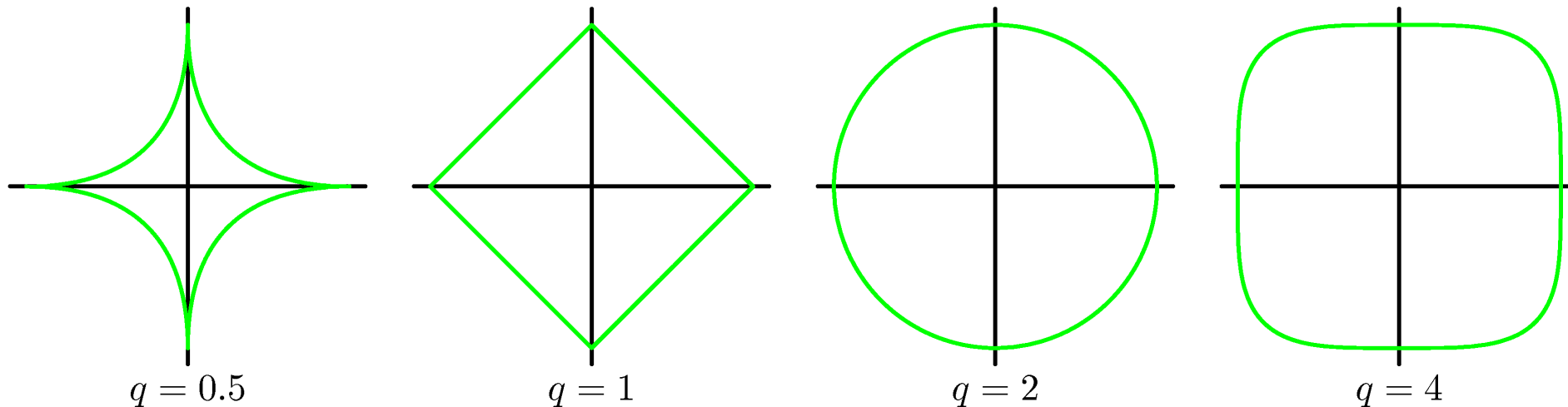
- Minimizing this yields

$$\mathbf{w} = \left( \lambda \mathbf{I} + \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{t}.$$

# Regularized Least Squares

- More generic regularizer form

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$



**Lasso**

# Multiple Outputs

- Analogously to the multiple output case we have:

$$\begin{aligned} \mathbf{y}(\mathbf{x}, \mathbf{w}) &= \mathbf{W}^T \phi(\mathbf{x}) & p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta) &= \mathcal{N}(\mathbf{t}|\mathbf{y}(\mathbf{W}, \mathbf{x}), \beta^{-1} \mathbf{I}) \\ & & &= \mathcal{N}(\mathbf{t}|\mathbf{W}^T \phi(\mathbf{x}), \beta^{-1} \mathbf{I}). \end{aligned}$$

- Given observed inputs,  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , and targets,  $\mathbf{T} = [\mathbf{t}_1, \dots, \mathbf{t}_N]^T$  we obtain the log likelihood function

$$\begin{aligned} \ln p(\mathbf{T}|\mathbf{X}, \mathbf{W}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(\mathbf{t}_n | \mathbf{W}^T \phi(\mathbf{x}_n), \beta^{-1} \mathbf{I}) \\ &= \frac{NK}{2} \ln \left( \frac{\beta}{2\pi} \right) - \frac{\beta}{2} \sum_{n=1}^N \|\mathbf{t}_n - \mathbf{W}^T \phi(\mathbf{x}_n)\|^2. \end{aligned}$$

# Multiple Outputs

- Maximizing with respect to  $W$ , we obtain

$$\mathbf{W}_{\text{ML}} = \left( \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{T}.$$

- If we consider a single target variable,  $t_k$ , we see that

$$\mathbf{w}_k = \left( \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{t}_k = \Phi^\dagger \mathbf{t}_k$$

- where  $\mathbf{t}_k = [t_{1k}, \dots, t_{Nk}]^T$ , which is identical with the single output case.

# Bias- Variance

- So far in linear models for regression: We have assumed that the form and number of basis functions are both fixed

***Limiting the number of basis functions in order to avoid over-fitting***

*V<sub>s</sub>*

***limiting the flexibility of the model to capture interesting and important trends in the data***

- The introduction of regularization terms can control over-fitting for models with many parameters
- New question: Suitable values of these parameters



# Bias- Variance

- Recall the *expected squared loss*,

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \underbrace{\int \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt}_{\text{noise}}$$

- where

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int t p(t|\mathbf{x}) dt.$$

- The second term of  $\mathbb{E}[L]$  corresponds to the noise inherent in the random variable  $t$ .
- What about the first term?

# Bayesian Linear Regression

- Define a conjugate prior over  $\mathbf{w}$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0).$$

- Combining this with the likelihood function and using results for marginal and conditional Gaussian distributions, gives the posterior

$$p(\mathbf{w} | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$$

- where

$$\mathbf{m}_N = \mathbf{S}_N \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t} \right)$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi.$$

# Bayesian Linear Regression

- A common choice for the prior is (zero-mean isotropic Gaussian)

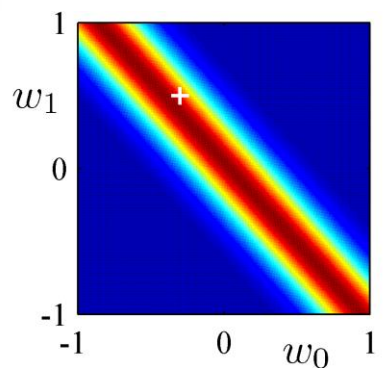
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$$

for which

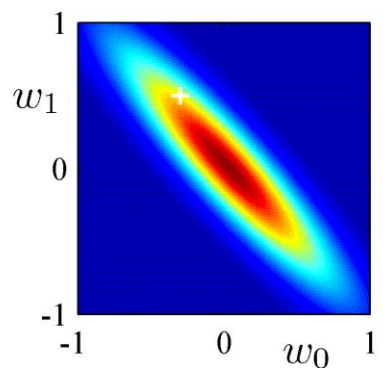
$$\begin{aligned} \mathbf{m}_N &= \beta \mathbf{S}_N \Phi^T \mathbf{t} \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \Phi^T \Phi. \end{aligned}$$

- Next we consider an example ...

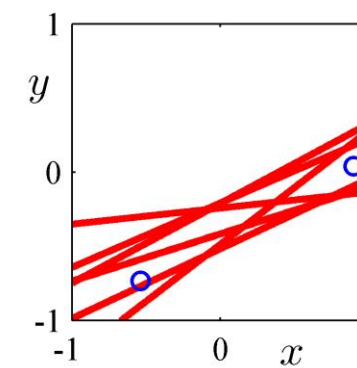
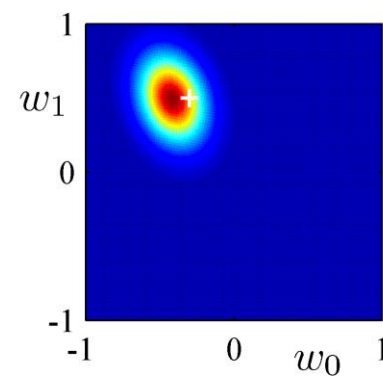
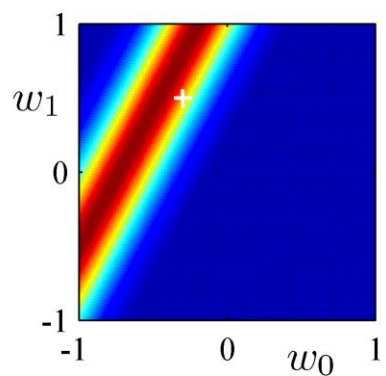
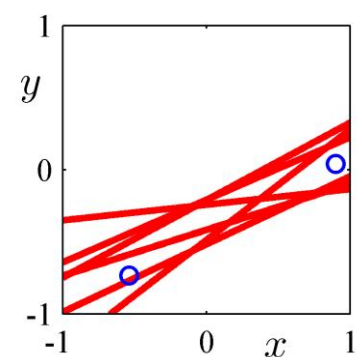
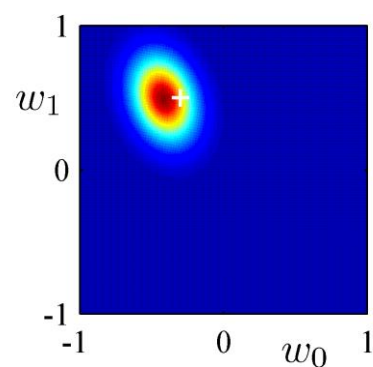
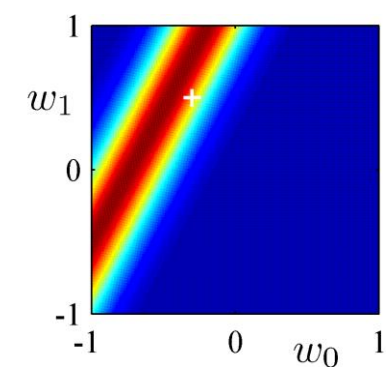
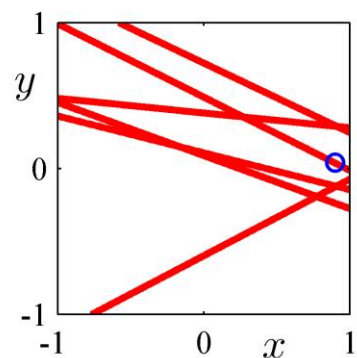
Likelihood



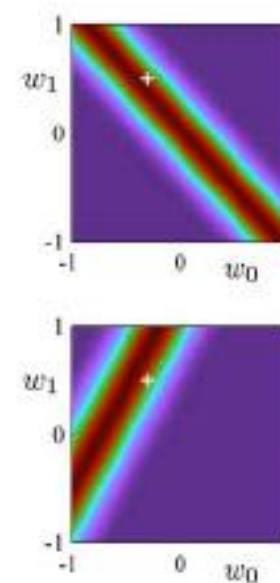
Posterior



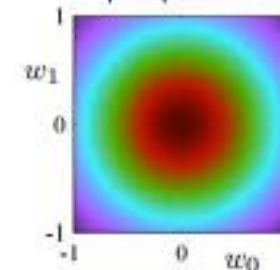
Data Space



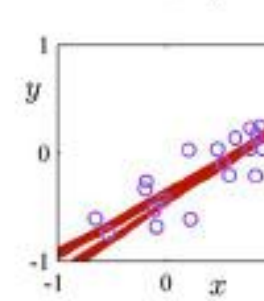
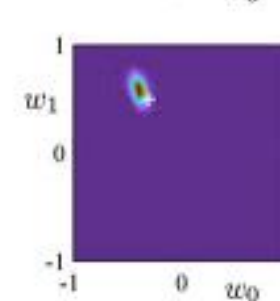
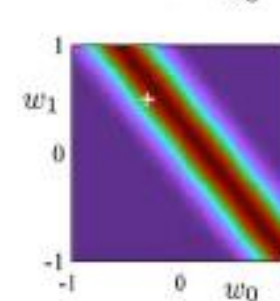
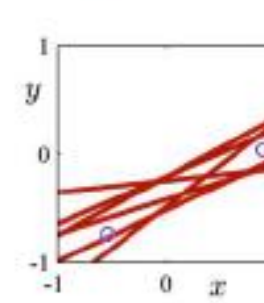
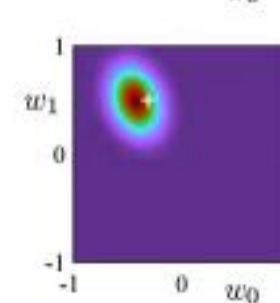
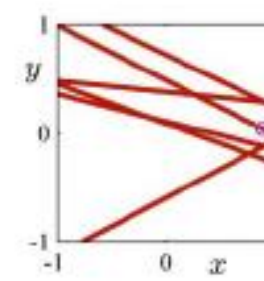
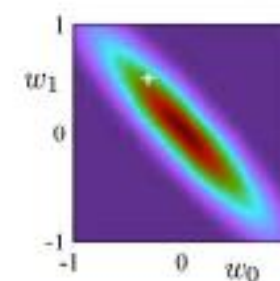
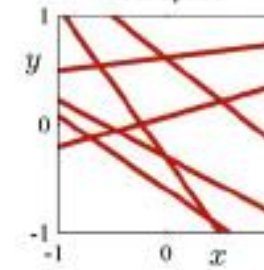
likelihood



prior/posterior



data space



# Thank You

