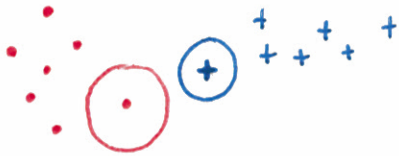
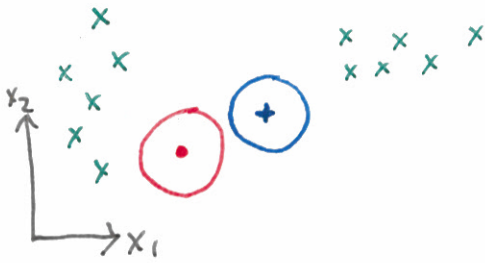


# EM for Mixtures of Gaussians



0. Initialize parameters

$$\theta = \{ \pi, \{ \mu^{(k)}, \Sigma^{(k)} \} \}$$

1. E-Step

Fix soft responsibilities:

$$\Gamma_k^{(n)} = P(z^{(n)} = k \mid \underline{x}^{(n)}, \theta)$$

2. M-Step

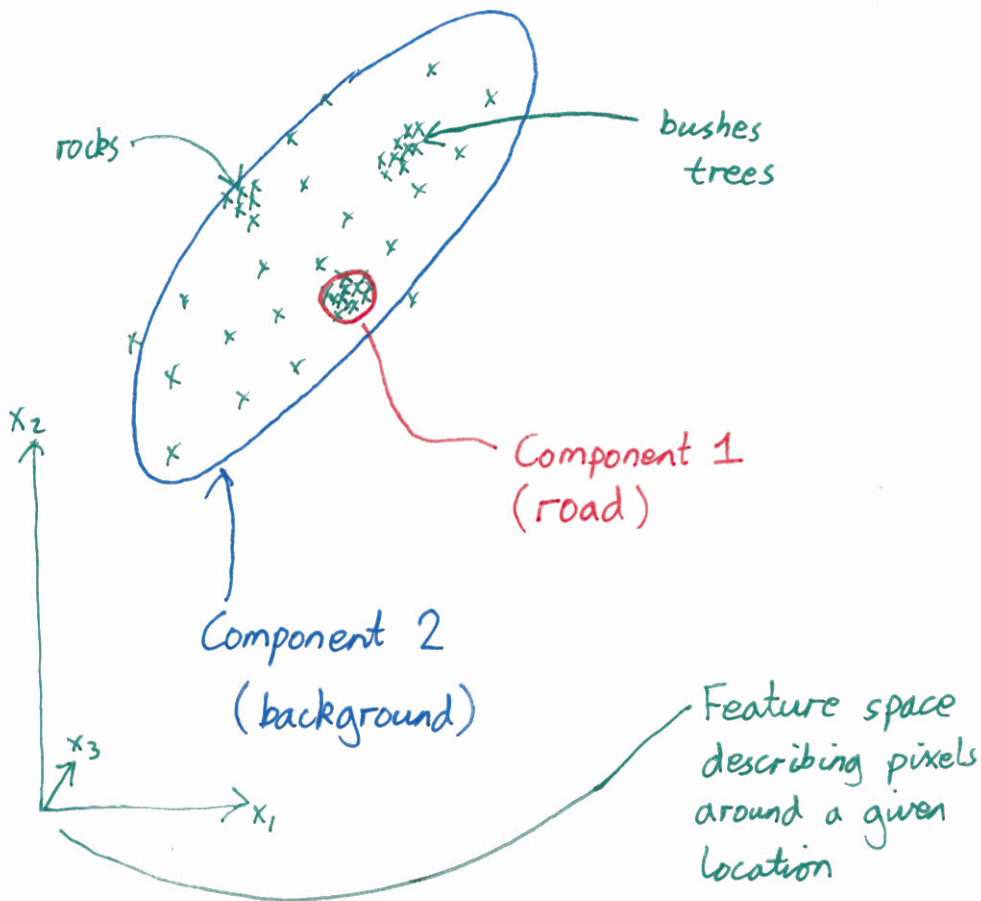
Fit  $\theta$  using  $\{ \Gamma_k^{(n)} \}$

$$\Gamma_k = \sum_n \Gamma_k^{(n)}, \quad \pi_k = \frac{\Gamma_k}{N}, \quad \underline{\mu}^{(k)} = \frac{1}{\Gamma_k} \sum_n \Gamma_k^{(n)} \underline{x}^{(n)}$$

3. Go to 1. or stop.

$$\Sigma^{(k)} = \frac{1}{\Gamma_k} \sum_n \Gamma_k^{(n)} (\underline{x}^{(n)} - \underline{\mu}^{(k)}) (\underline{x}^{(n)} - \underline{\mu}^{(k)})^T$$

# Use in "Stanley" car



## Interpretation

We are using an approx. posterior

$$Q(z^{(n)} = k) = r_k^{(n)} = P(z^{(n)} = k | \underline{x}^{(n)}, \theta^{(old)})$$

Only correct when  $\theta = \theta^{(old)}$

## Compare distributions

$$D_{KL}(Q \parallel P) = \sum_{\underline{z}} Q(\underline{z}) \log \frac{Q(\underline{z})}{P(\underline{z} | X, \theta)} \geq 0$$

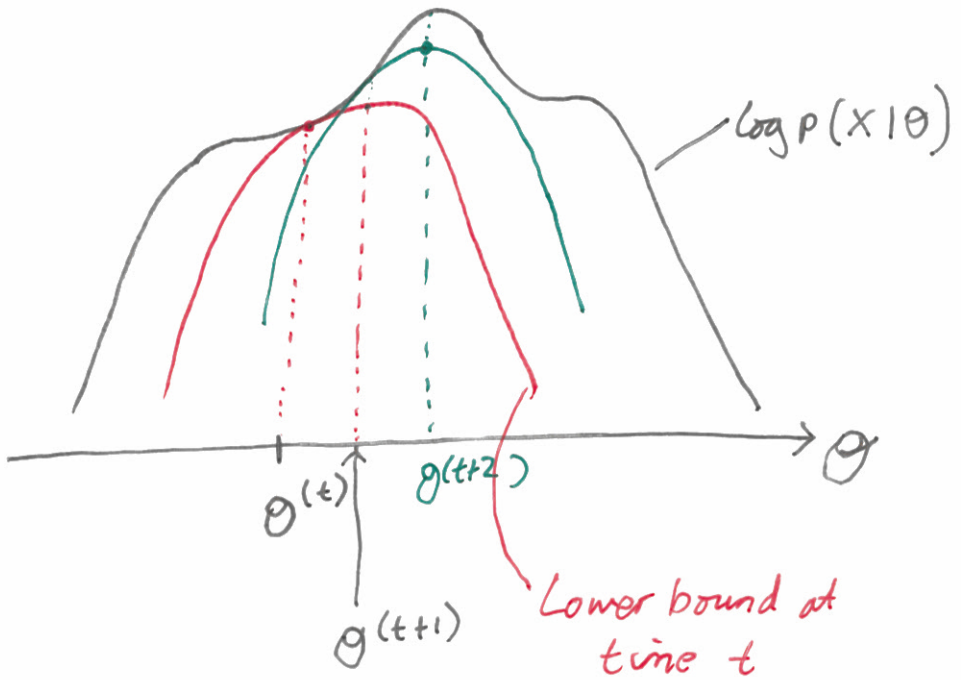
$$\text{Use } P(\underline{z} | X, \theta) = \frac{P(\underline{z}, X | \theta)}{P(X | \theta)} \leftarrow \text{Likelihood}$$

$$\Rightarrow \sum_{\underline{z}} Q(\underline{z}) \log \frac{Q(\underline{z})}{P(X, \underline{z} | \theta)} \geq -\log P(X | \theta)$$

$\Rightarrow$  Bound on log-likelihood

Tight if  $\theta = \theta^{(old)}$  used to set  $Q$

# Bound-based optimizer



## Newton's Method

Cost function  $E(\underline{w})$

Gradients  $\underline{g} = \nabla_{\underline{w}} E(\underline{w})$

Hessian  $H_{ij} = \frac{\partial^2 E}{\partial w_i \partial w_j}$

Initialize  $\underline{w}^{(0)}$

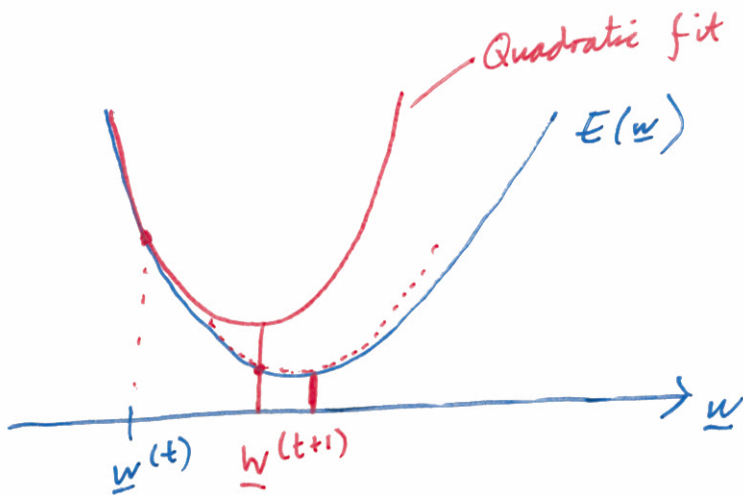
$$\underline{w}^{(t+1)} = \underline{w}^{(t)} - H^{-1} \underline{g}$$

If we have a quadratic cost:

$$E(\underline{w}) = \frac{1}{2} (\underline{w} - \underline{w}^{(*)})^T H (\underline{w} - \underline{w}^{(*)}) + \text{const.}$$

$$\text{Here } \underline{g} = H (\underline{w} - \underline{w}^{(*)})$$

$$\begin{aligned} \underline{w}^{(t+1)} &= \underline{w}^{(t)} - \cancel{H^{-1} H} (\underline{w}^{(t)} - \underline{w}^{(*)}) \\ &= \underline{w}^{(*)} \end{aligned}$$



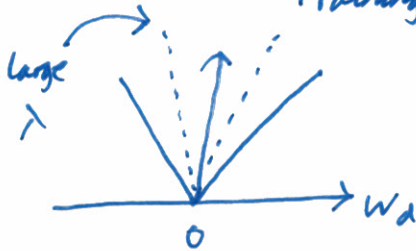
Why use other optimizers?

- Convergence?
- Not needing to tune...
- SGD can't give "sparse" solutions

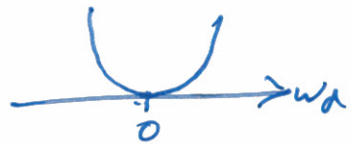
Some  $w_d = 0$

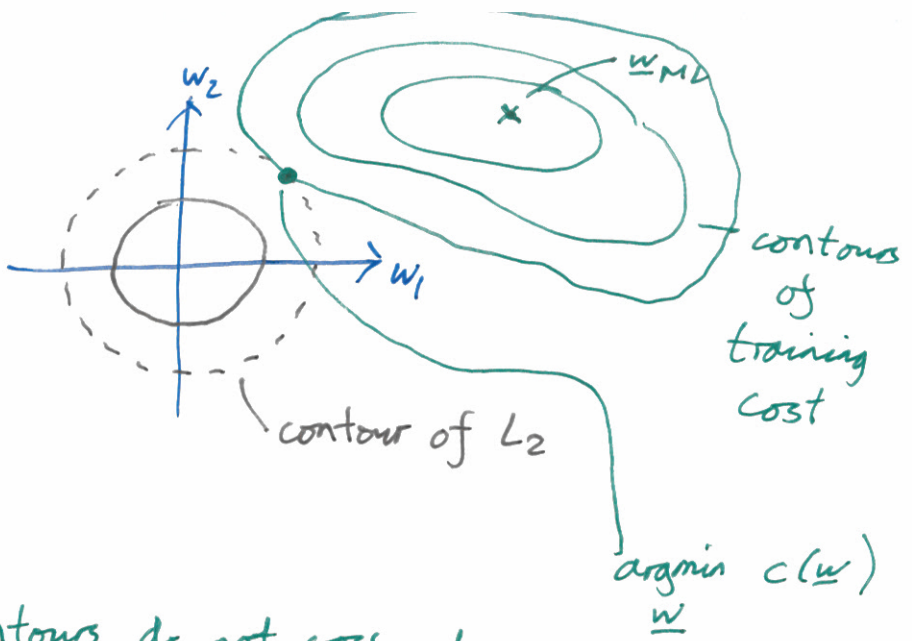
# L1 Regularization

$$c(\underline{w}) = \underbrace{E(\underline{w})}_{\text{Training error}} + \underbrace{\lambda \sum_d |w_d|}_{\lambda \|\underline{w}\|_1}$$

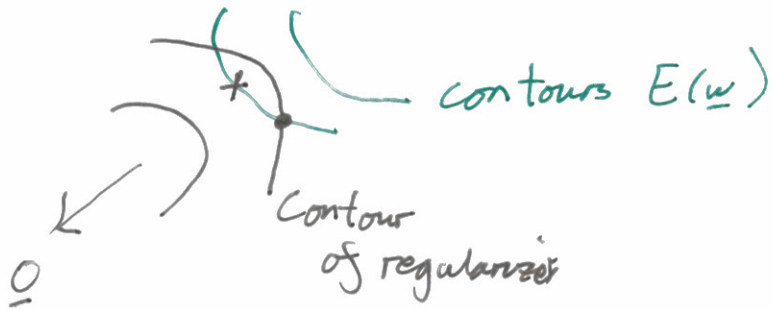


L2 Regularizer

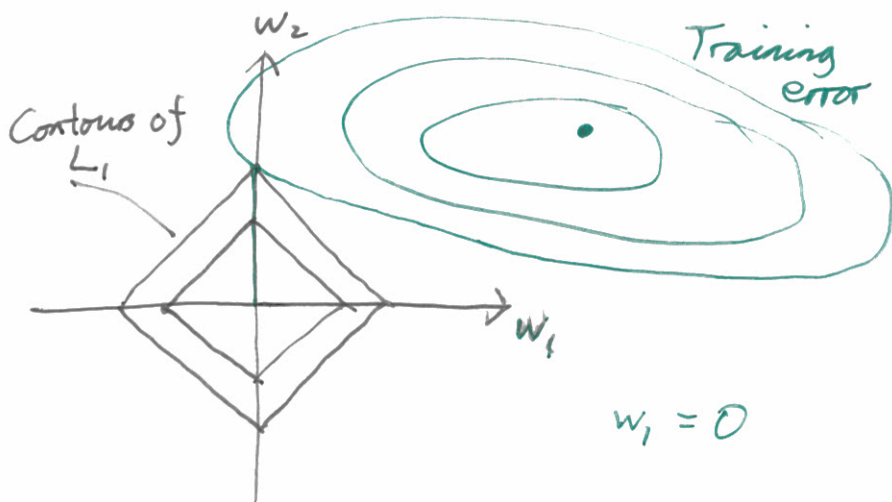




Contours do not cross at optimum:

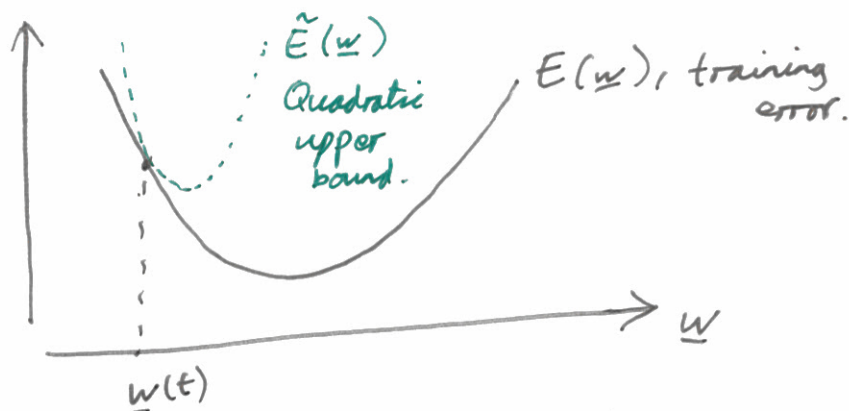






There are many ways to fit  $L_1$

One way are proximal methods:



$$\tilde{E}(\underline{w}) = \tilde{E}(\underline{w}) + \lambda \|\underline{w}\|_1$$

Can fit in closed form, some  $w_d = 0$ .

From Bayesian view, predictions are never  
sparse

$$P(y | \underline{x}, D) = \int p(y | \underline{x}, \underline{w}) \underbrace{p(\underline{w} | D)}_{> 0} d\underline{w}$$

$$P(w_a \neq 0 | D) > 0$$