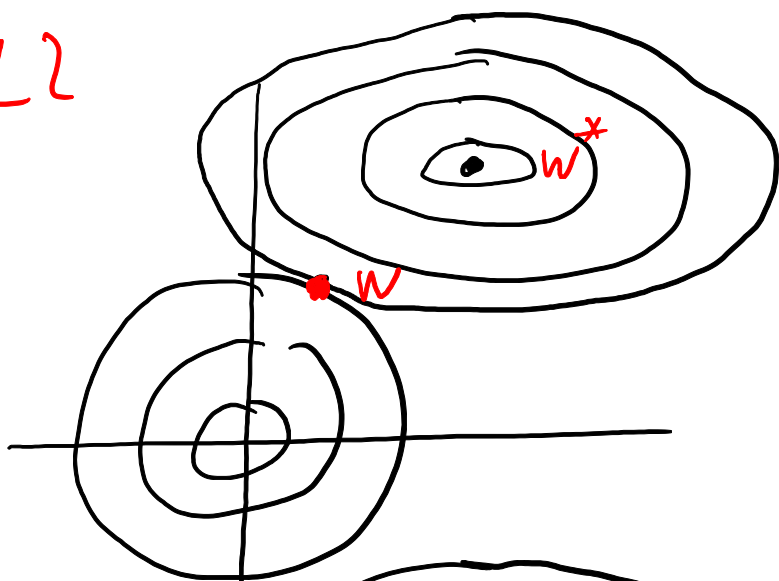- DL works well with high dimension
- Deeper network requires narrower layer
- Regularize bias will undertit model
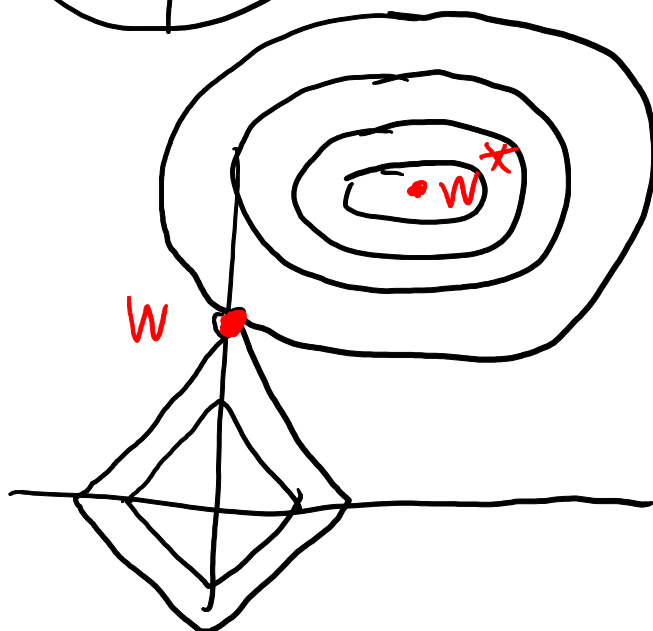
$$W_{n+1} = W_n - \varepsilon(\alpha W_n + \nabla_{W_n} J) \overset{L_2}{=} W_n - \varepsilon(\alpha\, sign(W_n) + \nabla_{W_n} J) \overset{L_1}{=}$$

$$= (1-\varepsilon\alpha)w - \varepsilon\nabla_{W_n} J = (1-\varepsilon\alpha)sign(W_n) - \varepsilon\nabla_{W_n} J$$

- L2



$w^*$ overfits

L1

- Minibatch SGD:

  10000 samples, 200 batches.

  SGD compute 200 more calculations,

  but reduce error by factor of $\sqrt{50}$

- Using GPU, batch size should be in size

  of $2^k$, ie 32, 64, 128, 256

- Batch is also considered as regularizer

- Hessian matrix requires more samples per

  batch, but size of $w$ can be too large!!!

  So we don't use 2nd order method.

- <span style="color:red">Read on momentum & initialization</span>

- Initializing weights should break symetry.

  $$W_0 \sim N(0, \sigma^2) \quad \& \quad U\left(-\frac{1}{\sqrt{m}}, \frac{1}{\sqrt{m}}\right) \& U\left(-\frac{1}{\sqrt{m+n}}, \frac{1}{\sqrt{m+n}}\right)$$

- Initializing bais to match marginal distribution

  of the data; bias should not be zero.