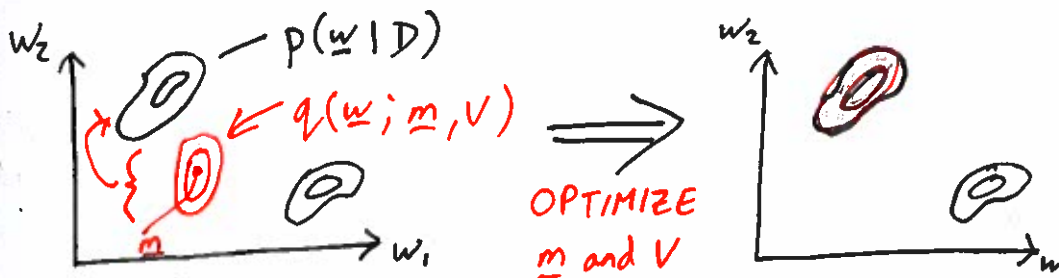


Variational Inference



Minimize $D_{KL}(q \parallel p(w|D))$

$$p(w|D) = \frac{p(D|w) p(w)}{p(D)}$$

$$D_{KL} = \mathbb{E}_q \left[\log \frac{q(w; m, V)}{p(w|D)} \right] \geq 0$$

$\mathbb{E}[\text{"energy"}]$

$$= \underbrace{-\mathbb{E}_q[\log p(D|w)]}_{\mathbb{E}[\text{neg. log likelihood}]} - \underbrace{\mathbb{E}_q[\log p(w)]}_{D_{KL}(q \parallel p(w))} + \underbrace{\mathbb{E}_q[\log q(w; m, V)]}_{-\text{Entropy}[q]} + \log p(D)$$

\uparrow
Marginal Likelihood

$$= \mathcal{J} + \log p(D) \geq 0 \quad (\text{Gibbs' inequality})$$

Marginal Likelihood bound

$$\log p(D) \geq -\mathcal{J}$$

Minimize \mathcal{J} wrt (m, V) and σ_w^2, \dots
Approx. Posterior well.

Find good model

Week 11 events

No MLPR Lectures

(last lecture Thurs 21 Nov.)

Ed-Intelligence have two
(Links also in "week 11" of mlpr notes) events!

Mini NeurIPS

6pm Wed 27 Nov, AT LT 2

tinyurl.com/mini-neurips-2019

Biases, failure + fairness in AI

6pm Fri 29 Nov, AT LT 5

to-err-is-machine.eventbrite.co.uk

Minimize J

Need tricks to make SGD work

Trick # 1, Reparameterize to make unconstrained

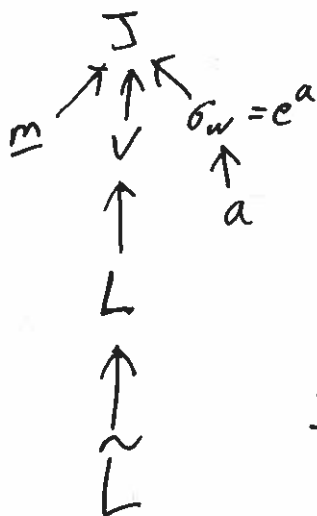
SGD on σ_w break, $\sigma_w < 0$ undefined

$\sigma_w \approx 0$ unstable

Set $\sigma_w = e^a$, optimize a

V positive definite, symmetric

Compute graph



$$V = LL^T$$

$$L_{ij} = \begin{cases} e^{\tilde{L}_{ii}} & i=j \\ \tilde{L}_{ij} & i>j \\ 0 & i<j \end{cases}$$

If could evaluate J ,
backprop, do SGD on
 m, \tilde{L}, a

Evaluating the cost, J

$D_{KL}(q \parallel p(\underline{w}))$ this is "easy"

↳ often Gaussian (Look it up in Matrix Cookbook)

Likelihood term

$$\begin{aligned} & \mathbb{E}_q[\log p(D|\underline{w})] \\ &= \mathbb{E}_q\left[\sum_{n=1}^N \log p(y^{(n)} | \underline{x}^{(n)}, \underline{w})\right] \end{aligned}$$

For logistic regression \rightarrow 1D integral do numerically

Trick # 2 "Reparameterization trick"

Stochastic estimate

$$\begin{aligned} & \mathbb{E}_{N(\underline{w}; \underline{m}, \underline{V})} [f(\underline{w})] \approx \frac{1}{S} \sum_s f(\underline{w}^{(s)}) \\ & \quad \underline{w}^{(s)} \sim N(\dots) \\ &= \mathbb{E}_{N(\underline{v}; \underline{0}, \underline{I})} [f(\underline{L}\underline{v} + \underline{m})] \\ &\approx f(\underline{L}\underline{v} + \underline{m}), \quad \underline{v} \sim N(\underline{0}, \underline{I}) \\ & \quad \text{(Using only sample!)} \end{aligned}$$

Estimate of gradients

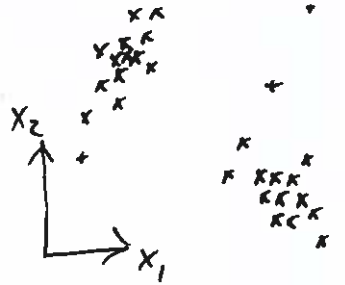
$$\begin{aligned}\nabla_{\underline{m}} \mathbb{E}_{N(\underline{w}; \underline{m}, \underline{V})} [f(\underline{w})] \\ \approx \nabla_{\underline{m}} f(\underline{L}\underline{v} + \underline{m}) = \nabla_{\underline{w}} f(\underline{w}) \Big|_{\underline{w} = \underline{L}\underline{v} + \underline{m}}\end{aligned}$$

$$\begin{aligned}\nabla_{\underline{L}} \mathbb{E}_{N(\underline{w}; \underline{m}, \underline{V})} [f(\underline{w})] \\ \approx \nabla_{\underline{L}} f(\underline{L}\underline{v} + \underline{m}) \\ = \nabla_{\underline{w}} f(\underline{w}) \Big|_{\underline{w} = \underline{L}\underline{v} + \underline{m}} \underline{v}^T\end{aligned}$$

Use same $\nabla_{\underline{w}} f(\underline{w})$ derivatives as normal but on noisy weights.

Mixtures of Gaussians

- Uses:
- Clustering
 - Model non-Gaussian distributions



Model:

Same as Gaussian Bayes classifier

Except "labels" are hidden or "latent" variables

Cluster choice $z^{(n)} \sim \text{Discrete}(\pi)$

$$z^{(n)} \in \{1, 2, 3, \dots, k\}$$

↑
+ve vector
sums to 1

Features

$$\underline{x}^{(n)} | z^{(n)} = k \sim \mathcal{N}(\underline{x}^{(n)}; \underline{\mu}^{(k)}, \Sigma^{(k)})$$

Likelihood

Given one observation: $p(\underline{x}^{(n)} | \theta = (\pi, \{\underline{\mu}^{(k)}, \Sigma^{(k)}\}))$

$$= \sum_{k=1}^k p(\underline{x}^{(n)}, z^{(n)} = k | \theta)$$

$$= \sum_k p(\underline{x}^{(n)} | z^{(n)} = k, \theta) p(z^{(n)} = k | \theta)$$

$$= \sum_k \mathcal{N}(\underline{x}^{(n)}; \underline{\mu}^{(k)}, \Sigma^{(k)}) \pi_k$$

Unsupervised learning, Clustering

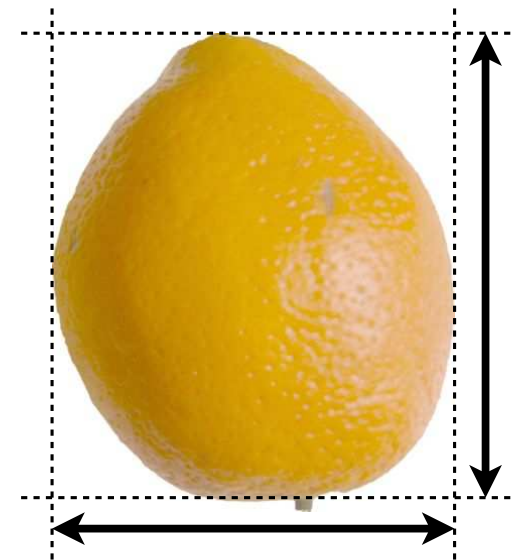
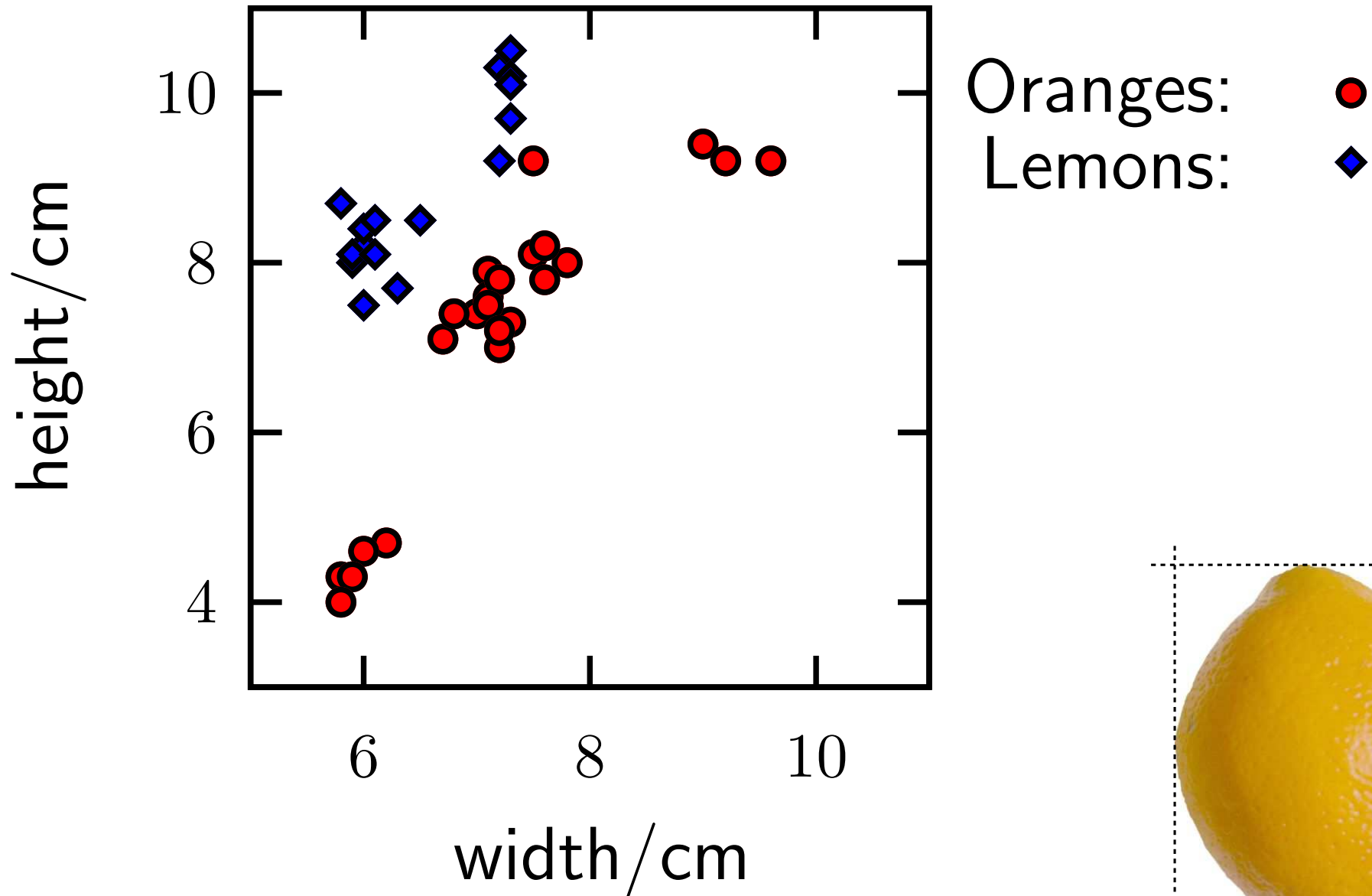
“Human brains are good at finding regularities in data.

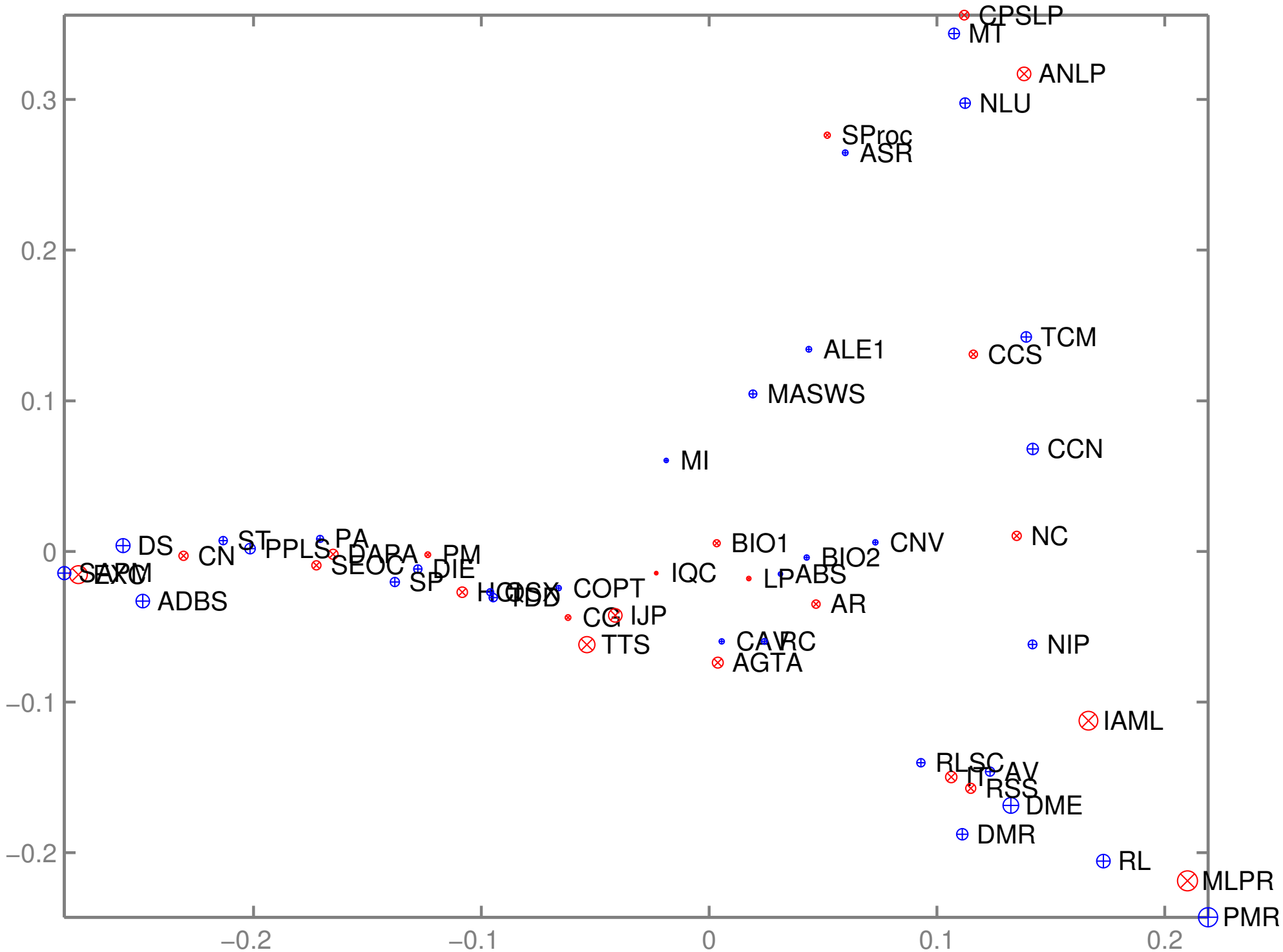
One way of expressing regularity is to put a set of objects into groups that are similar to each other.

For example, biologists have found that most objects in the natural world fall into one of two categories: things that are brown and run away, and things that are green and don't run away. The first group they call animals, and the second, plants.”

— David MacKay, ITILA textbook p284

Oranges and Lemons data





Stanley



Stanford Racing Team; DARPA 2005 challenge

<http://robots.stanford.edu/talks/stanley/>

How to stay on a road?

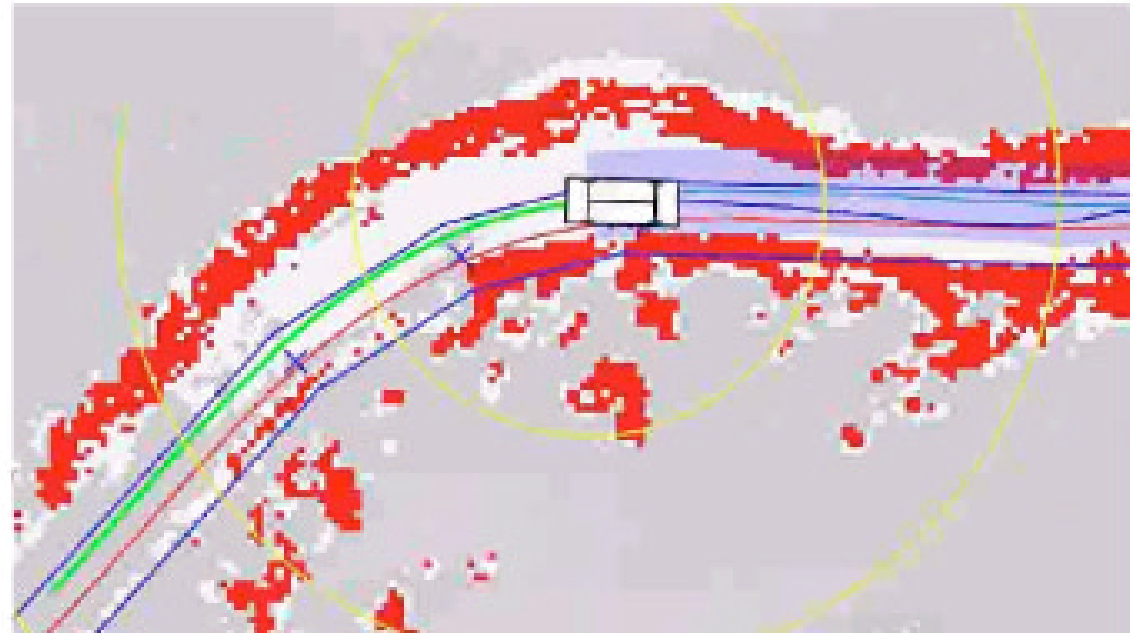


Perception and intelligence

(a) Beer Bottle Pass

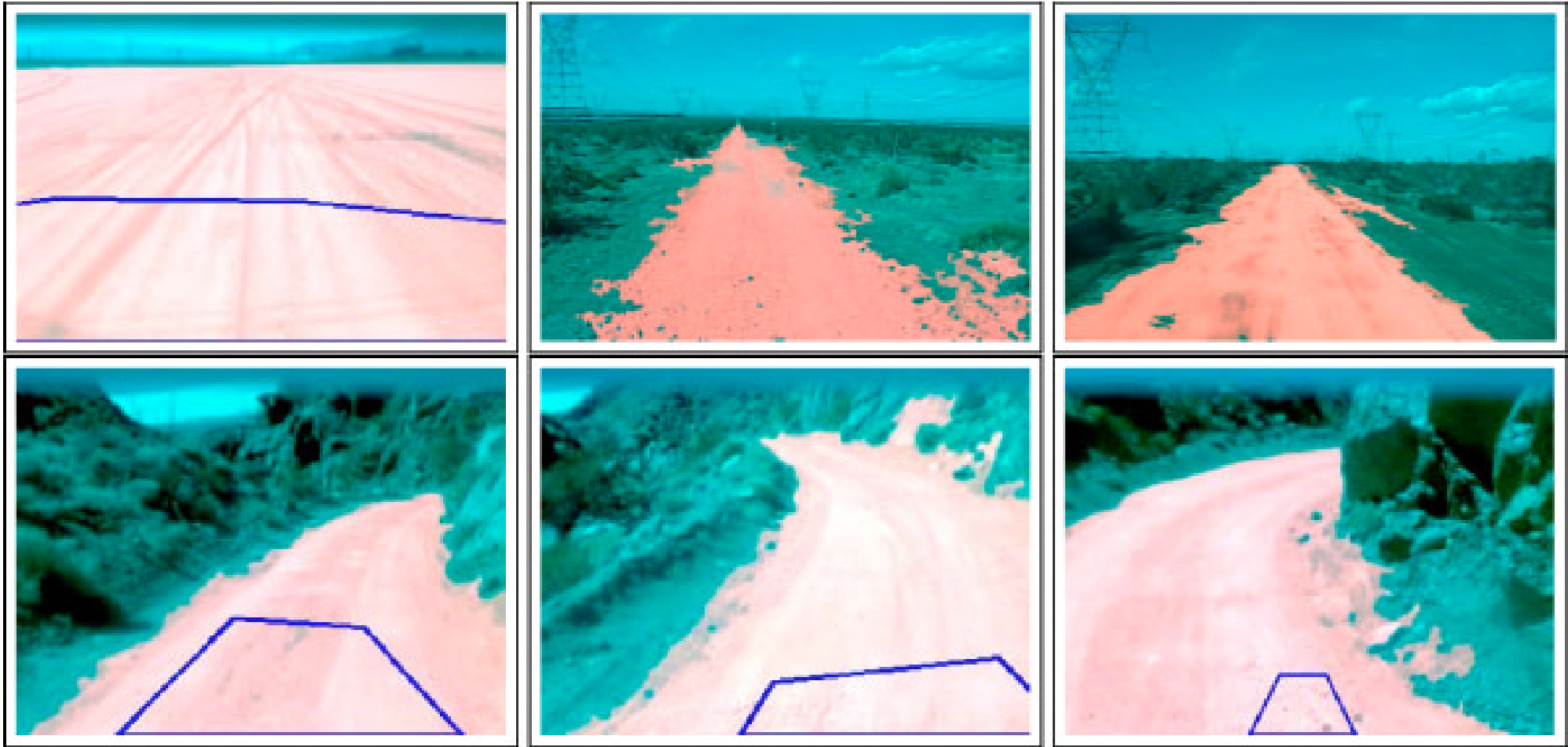


(b) Map and GPS corridor



It would look pretty stupid to run off the road, just because the trip planner said so.

Clustering to stay on the road



Stanley used a Gaussian mixture model.

The cluster just in front is road (unless we already failed).