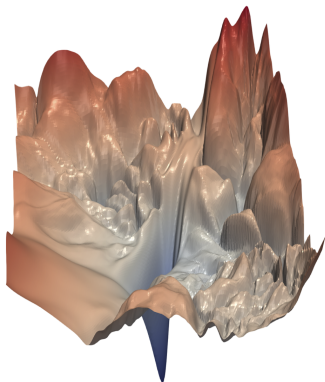


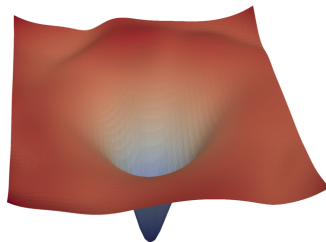
Outline for today

- ▶ Visualisation of some network loss landscapes
 - ▶ ResNet-56 with and without skip connection
 - ▶ VGG9 dependence on weight decrease and normalisation
 - ▶ ResNet depth dependence and skip connections
- ▶ Early theoretical results on over parameterisation and number of connected components of loss level curves.
- ▶ Removing the weight nonlinearity through convexification.

Loss landscape example: ResNet-56 (Li et al. 18'¹)



(a) without skip connections



(b) with skip connections

Figure 1: The loss surfaces of ResNet-56 with/without skip connections. The proposed filter normalization scheme is used to enable comparisons of sharpness/flatness between the two figures.

¹<http://papers.nips.cc/paper/7875-visualizing-the-loss-landscape-of-neural-nets.pdf>

Loss landscape example: VGG9 (Li et al. 18'²)

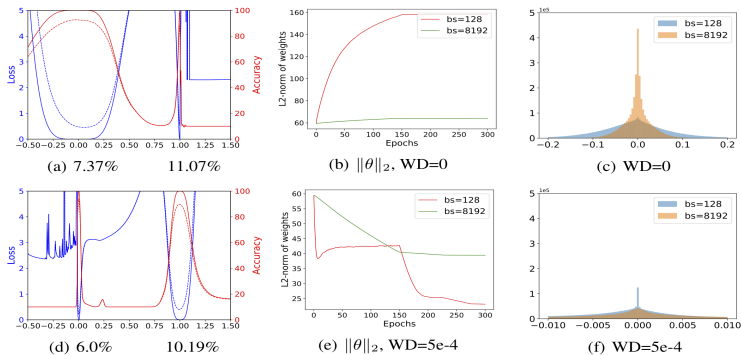


Figure 2: (a) and (d) are the 1D linear interpolation of VGG-9 solutions obtained by small-batch and large-batch training methods. The blue lines are loss values and the red lines are accuracies. The solid lines are training curves and the dashed lines are for testing. Small batch is at abscissa 0, and large batch is at abscissa 1. The corresponding test errors are shown below. (b) and (e) shows the change of weights norm $\|\theta\|_2$ during training. When weight decay is disabled, the weight norm grows steadily during training without constraints (c) and (f) are the weight histograms, which verify that small-batch methods produce more large weights with zero weight decay and more small weights with non-zero weight decay.

²<http://papers.nips.cc/paper/7875-visualizing-the-loss-landscape-of-neural-nets.pdf>

Loss landscape example: VGG9 normalized (Li et al. 18³)

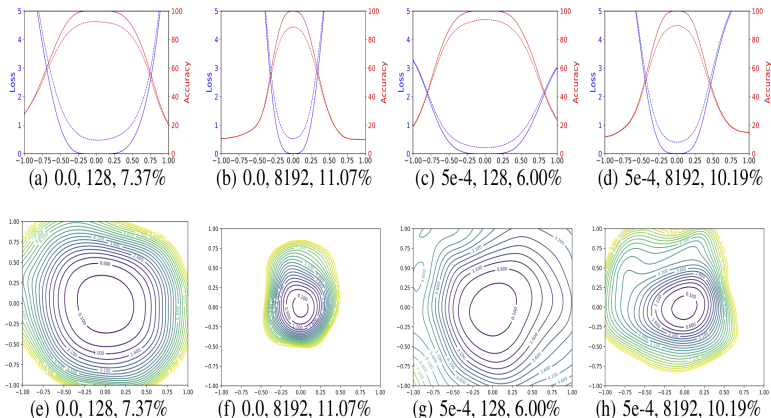


Figure 3: The 1D and 2D visualization of solutions obtained using SGD with different weight decay and batch size. The title of each subfigure contains the weight decay, batch size, and test error.

³<http://papers.nips.cc/paper/>

7875-visualizing-the-loss-landscape-of-neural-nets.pdf

Loss landscape example: ResNet skip (Li et al. 18'⁴)

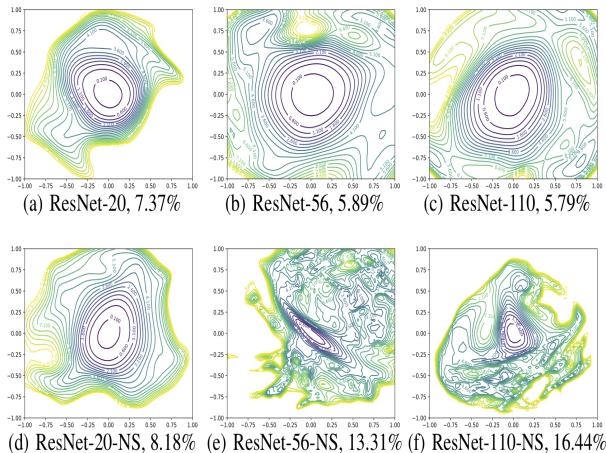


Figure 5: 2D visualization of the loss surface of ResNet and ResNet-noshort with different depth.

⁴<http://papers.nips.cc/paper/7875-visualizing-the-loss-landscape-of-neural-nets.pdf>

Loss landscape example: ResNet width (Li et al. 18⁵)

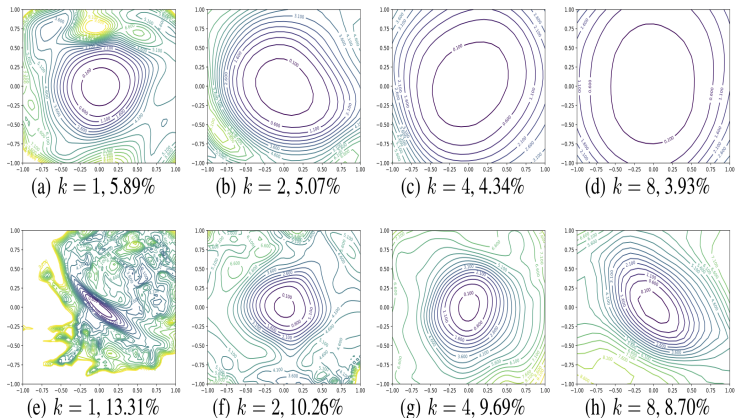


Figure 6: Wide-ResNet-56 on CIFAR-10 both with shortcut connections (top) and without (bottom). The label $k = 2$ means twice as many filters per layer. Test error is reported below each figure.

⁵<http://papers.nips.cc/paper/>

7875-visualizing-the-loss-landscape-of-neural-nets.pdf

Topology of loss landscape (Freeman et al. 16'⁶)

Consider our loss function: $\mathcal{L}(\theta; X, Y) = n^{-1} \sum_{\mu=1}^n l(\theta; x_{\mu}, y_{\mu})$
and its associated level set

$$\Omega_{\mathcal{L}}(\lambda) = \{\theta : \mathcal{L}(\theta; X, Y) \leq \lambda\}$$

Of particular interest are the number of connected components, say N_{λ} , in $\Omega_{\mathcal{L}}(\lambda)$. If $N_{\lambda} = 1$ for all λ then $\mathcal{L}(\theta; X, Y)$ has no isolated local minima and any descent method can obtain a global minima.

If $N_{\lambda} > 1$ there may be “spurious valleys” in which the minima in the connected component does not achieve the global minima.

⁶<https://arxiv.org/pdf/1611.01540.pdf>

Topology of loss landscape (Freeman et al. 16⁷)

Linear network: single component

Let $H(x; \theta)$ be an L layer net given by $h^{(\ell)} = W^{(\ell)}h^{(\ell-1)}$ with $W^{(\ell)} \in \mathbb{R}^{n_\ell \times n_{\ell-1}}$, then if $n_\ell > \min(n_0, n_L)$ for $0 < \ell < L$, the sum of squares loss function has a single connected component

ReLU network: multiple components

Let $H(x; \theta)$ be an L layer net given by $h^{(\ell)} = \sigma(W^{(\ell)}h^{(\ell-1)})$ with $W^{(\ell)} \in \mathbb{R}^{n_\ell \times n_{\ell-1}}$ and $\sigma(\cdot) = \max(0, \cdot)$, then for any choice of n_ℓ there is a distribution of data (X, Y) such that there are more than one single connected component.

⁷<https://arxiv.org/pdf/1611.01540.pdf>

Topology of loss landscape: quadratic activation (Venturi et al. 16'⁸)

ReLU activation network: nearly connected

Consider a 2 layer ReLU network $H(x, \theta) = W^{(2)}\sigma(W^{(1)}x)$ with $W^{(1)} \in \mathbb{R}^{m \times n}$ and $W^{(2)} \in \mathbb{R}^m$, then for any two parameters θ_1 and θ_2 with $\mathcal{L}(\theta_i) \leq \lambda$ for $i = 1, 2$, then there is a path $\gamma(t)$ between θ_1 and θ_2 such that $\mathcal{L}(\theta_{\gamma(t)}) \leq \max(\lambda, m^{-1/n})$.

quadratic activation network: single component

Let $H(x, \theta)$ be an L layer net given by $h^{(\ell)} = \sigma(W^{(\ell)}h^{(\ell-1)})$ with $W^{(\ell)} \in \mathbb{R}^{n_\ell \times n_{\ell-1}}$ and quadratic activation $\sigma(z) = z^2$, then once the number of parameters $n_\ell \geq 3N^{2^\ell}$ where N is the number of data entries, then the sum of squares loss function has a single connected component. For the two layer case with a single quadratic activation this simplifies to $n > 2N$.

⁸<https://arxiv.org/pdf/1802.06384.pdf>

Convexifying the parameters pt. 1 (Zhang et al. 16'⁹)

Consider a two layer convolutional neural network composed of one convolutional layer followed by a fully connected layer.

Rather than working with x directly, form P vectors $z_p(x)$ for $p = 1, \dots, P$ where $z_p(x)$ is the portion of x on patch p of the convolutional layer. Then the k^{th} component of $H(x, \theta)$ is given by

$$H(x, \theta)_k = \sum_{j=1}^r \sum_{p=1}^P \alpha_{k,j,p} \sigma(w_j^T z_p(x)).$$

Alternatively if we exclude the nonlinearity we can express this sum by

$$\sum_{j=1}^r \sum_{p=1}^P \alpha_{k,j,p} \sigma(w_j^T z_p(x)) = \sum_{j=1}^r Z(x) w_j$$

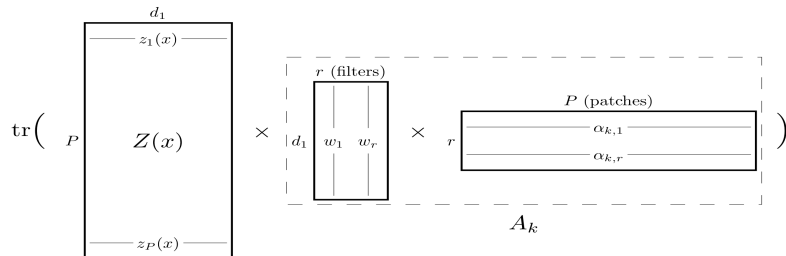
where $Z(x)$ has $z_p(x)$ as its p^{th} row.

⁹<https://arxiv.org/pdf/1609.01000.pdf>

Convexifying the parameters pt. 2 (Zhang et al. 16'¹⁰)

Using the trace formula this can be further condensed to

$$H(x, \theta)_k = \text{tr} \left(Z(x) \left(\sum_{j=1}^r w_j \alpha_{k,j}^T \right) \right) = \text{tr} (Z(x) A_k)$$



The network parameters are given by A_k nonlinearity is imposed by the A_k having rank r , and we can express all of the parameters of the matrix by A which is similarly rank r .

¹⁰<https://arxiv.org/pdf/1609.01000.pdf>

Convexifying the parameters pt. 3 (Zhang et al. 16'¹¹)

One can impose the network structure through A , but remove the non-convex rank constraint by replacing a convexification, that is the sum of the singular values of A (Schatten-1, or nuclear, norm).

If the convolutional filters and fully connected rows are uniformly bounded in ℓ^2 by B_1 and B_2 respectively, then one can replace then the sum of the singular values of A are bounded by $B_1 B_2 r \sqrt{n}$ where n is the network output dimension and the network parameters can be considered by varying the nuclear norm bound between 0 and $B_1 B_2 r \sqrt{n}$.

The resulting learning programme is fully convex and can be efficiently solved. The above can be extended to nonlinear activations and multiple layers, learning one layer at a time.

¹¹<https://arxiv.org/pdf/1609.01000.pdf>

Convexified CNN: MNIST (Zhang et al. 16'¹²)

	basic	rand	rot	img	img+rot
SVM _{r_{bf}} [44]	3.03%	14.58%	11.11%	22.61%	55.18%
NN-1 [44]	4.69%	20.04%	18.11%	27.41%	62.16%
CNN-1 (ReLU)	3.37%	9.83%	18.84%	14.23%	45.96%
CCNN-1	2.38%	7.45%	13.39%	10.40%	42.28%
TIRBM [38]	-	-	4.20%	-	35.50%
SDAE-3 [44]	2.84%	10.30%	9.53%	16.68%	43.76%
ScatNet-2 [8]	1.27%	12.30%	7.48%	18.40%	50.48%
PCANet-2 [9]	1.06%	6.19%	7.37%	10.95%	35.48%
CNN-2 (ReLU)	2.11%	5.64%	8.27%	10.17%	32.43%
CNN-2 (Quad)	1.75%	5.30%	8.83%	11.60%	36.90%
CCNN-2	1.38%	4.32%	6.98%	7.46%	30.23%

Table 1: Classification error on the basic MNIST and its four variations. The best performance within each block is bolded. The tag “ReLU” and “Quad” means ReLU activation and quadratic activation, respectively.

¹²<https://arxiv.org/pdf/1609.01000.pdf>

Convexified CNN: CIFAR10 (Zhang et al. 16'¹³)

	Error rate
CNN-1	34.14%
CCNN-1	23.62%
CNN-2	24.98%
CCNN-2	20.52%
SVM _{Fastfood} [27]	36.90%
PCANet-2 [9]	22.86%
CKN [30]	21.70%
CNN-3	21.48%
CCNN-3	19.56%

Table 3: Classification error on the CIFAR-10 dataset. The best performance within each block is bolded.

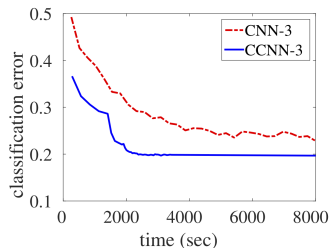


Figure 4: The convergence of CNN-3 and CCNN-3 on the CIFAR-10 dataset.

	CNN-1	CNN-2	CNN-3
Original	34.14%	24.98%	21.48%
Convexified	23.62%	21.88%	18.18%

Table 4: Comparing the original CNN and the one whose top convolution layer is convexified by CCNN. The classification errors are reported on CIFAR-10.

¹³<https://arxiv.org/pdf/1609.01000.pdf>